

Fair Infinitesimal Jackknife: Mitigating the influence of biased training data points without refitting

Neurips 2022

SeHyun Park

January 30, 2023

Seoul National University

- 1 Introduction
- 2 Preliminaries
- 3 Fair Classification through Post-Hoc Interventions
- 4 Experiment

Introduction

- ▶ Bias mitigation algorithm can be categorized into pre-processing, in-processing, and post-processing approaches.
 - ▶ Pre-processing and in-processing often require retraining a model from scratch and can be intractable in many real-world situation.
- ⇒ Post-processing approaches are the only viable option in such cases.

▶ Contribution:

- 1) Developing a post-processing fairness algorithm that improves the fairness characteristics of a pre-trained model without requiring it to be refit.
- 2) Developing **IHVP-WoodFisher**, a **WoodFisher** based on Inverse-Hessian Vector Product(IHVP) scheme for computing the fairness influence score.

Preliminaries

- ▶ Data set : $\mathcal{D} = \{\mathbf{z}_n = (\mathbf{x}_n, s_n, y_n)\}_{n=1}^N$
 - $\mathbf{x}_n \in \mathbb{R}^p$: feature
 - $s_n \in [0, 1, \dots, k]$: sensitive attribute
 - $y_n \in \mathcal{Y}$: response
- ▶ Parameter : $\boldsymbol{\theta} \in \Theta \subseteq \mathbb{R}^D$
- ▶ Model : $h_{\boldsymbol{\theta}}(\mathbf{x}_n) \in \mathcal{Y}$
- ▶ loss function $\ell : \Theta \times \mathcal{Y} \rightarrow \mathbb{R}$

► Weighted risk minimization problem

- Let $\mathbf{w} = [w_1, w_2, \dots, w_N]^T \in \mathbb{R}^N$ be weights vector.
- Then $\hat{\boldsymbol{\theta}}(\mathbf{w}) = \underset{\boldsymbol{\theta} \in \Theta}{\operatorname{argmin}} \frac{1}{N} \sum_{n=1}^N w_n \ell(h_{\boldsymbol{\theta}}(x_n), y_n)$
- Setting all the weights to one, $\mathbf{1} \stackrel{\text{def}}{=} [w_1 = 1, w_2 = 1, \dots, w_N = 1]^T$, then $\hat{\boldsymbol{\theta}}(\mathbf{1}) = \hat{\boldsymbol{\theta}} = \underset{\boldsymbol{\theta} \in \Theta}{\operatorname{argmin}} \frac{1}{N} \sum_{n=1}^N \ell(h_{\boldsymbol{\theta}}(x_n), y_n)$.
It recovers the standard empirical risk minimization problem.

- ▶ Since $\hat{\boldsymbol{\theta}}(\mathbf{w})$ is a function of the weights \mathbf{w} , we can form a first-order Taylor approximation to it about $\mathbf{1}$:

$$\hat{\boldsymbol{\theta}}(\mathbf{w}) = \hat{\boldsymbol{\theta}} + \nabla_{\mathbf{w}} \hat{\boldsymbol{\theta}}(\mathbf{w}) \Big|_{\mathbf{w}=\mathbf{1}} (\mathbf{w} - \mathbf{1}) + \mathcal{O} \left((\mathbf{w} - \mathbf{1})^2 \right)$$

- ▶ This first order Taylor approximation is often referred to as the **Infinitesimal Jackknife Approximation**.

- ▶ When $\hat{\theta}$ is a stationary point of $\frac{1}{N} \sum_{n=1}^N \ell(h_{\theta}(x_n), y_n) \stackrel{\text{let}}{=} L(\theta)$,

$$\left. \frac{d\hat{\theta}(\mathbf{w})}{dw_n} \right|_{\mathbf{w}=1} = -H^{-1}g_n$$

where $H \stackrel{\text{def}}{=} \nabla_{\theta}^2 L(\theta) \Big|_{\theta=\hat{\theta}}$, and $g_n \stackrel{\text{def}}{=} \nabla_{\theta} \ell(y_n, h_{\theta}(x_n)) \Big|_{\theta=\hat{\theta}}$

- ▶ To measure the influence of training instance on a differentiable functional, $M(\hat{\theta}(\mathbf{w}), \mathbf{w})$, of $\hat{\theta}(\mathbf{w})$, apply chain rule to arrive at,

$$\mathcal{I}_{M,n} \stackrel{\text{def}}{=} \left. \frac{dM(\hat{\theta}(\mathbf{w}), \mathbf{w})}{dw_n} \right|_{\mathbf{w}=1} = -\nabla_{\hat{\theta}} M(\hat{\theta}(\mathbf{w}), \mathbf{w}) \Big|_{\mathbf{w}=1}^T H^{-1}g_n$$

Fair Classification through Post-Hoc Interventions

- ▶ Develop : Post-processing fairness algorithm
- ▶ Given :
 - (1) a pre-trained model $\widehat{\theta}_{pre}$
 - (2) access to the training data and a validation set
 - (3) a twice differentiable loss function and a once differentiable surrogate to the fairness metric
 - (4) an invertible Hessian at a local optimum of the loss

► Two common fairness metrics

1. Demographic parity(DP)

$$- \Delta DP(\theta) = |P(h_{\theta}(X) = 1 | S = 1) - P(h_{\theta}(X) = 1 | S = 0)|$$

2. Equality of odds(EO)

$$- \Delta EO(\theta) = \sum_y |P(h_{\theta}(X) = 1 | S = 1, Y = y) - P(h_{\theta}(X) = 1 | S = 0, Y = y)|$$

► Smooth surrogate to ΔDP and ΔEO

$$- M_D^{\Delta DP}(\theta) = \left| \mathbb{E}_{p_{\mathcal{D}}(X=x|S=1)} [h_{\theta}(x)] - \mathbb{E}_{p_{\mathcal{D}}(X=x|S=0)} [h_{\theta}(x)] \right|$$

$$- M_D^{\Delta EO}(\theta) = \sum_y \left| \mathbb{E}_{p_{\mathcal{D}}(X=x|S=0, Y=y)} [h_{\theta}(x)] - \mathbb{E}_{p_{\mathcal{D}}(X=x|S=1, Y=y)} [h_{\theta}(x)] \right|$$

- ▶ Assume : use a held-out validation set $\mathcal{D}_{\text{val}} = \{\mathbf{x}_n, s_n, y_n\}_{n=1}^{N_{\text{val}}}$
- ▶ Influence function for group fairness :
 - on ΔDP

$$\mathcal{I}_{\Delta\text{DP},n} = -\nabla_{\hat{\boldsymbol{\theta}}} M_{\mathcal{D}_{\text{val}}}^{\Delta\text{DP}}(\hat{\boldsymbol{\theta}})^T H^{-1} g_n$$

- on ΔEO

$$\mathcal{I}_{\Delta\text{EO},n} = -\nabla_{\hat{\boldsymbol{\theta}}} M_{\mathcal{D}_{\text{val}}}^{\Delta\text{EO}}(\hat{\boldsymbol{\theta}})^T H^{-1} g_n$$

- ▶ Post-hoc improvement of pre-trained $\hat{\theta}$

$$\begin{aligned}\hat{\theta}_{\text{fair}} &\stackrel{\text{def}}{=} \hat{\theta}(\mathbf{w}_{\text{fair}}) = \hat{\theta} + \sum_{n=1}^N \left. \frac{d\theta(\mathbf{w})}{dw_n} \right|_{\mathbf{w}=\mathbf{1}} (w_n^{\text{fair}} - 1), \\ &= \hat{\theta} - \sum_{n=1}^N H^{-1} g_n (w_n^{\text{fair}} - 1)\end{aligned}$$

- ▶ By searching for a weight vector $\mathbf{w}_{\text{fair}} = [w_1^{\text{fair}}, w_2^{\text{fair}}, \dots, w_N^{\text{fair}}]^T \in \mathbb{R}^N$ such that optimizing $M_{\mathcal{D}_{\text{val}}}^b(\hat{\theta}(\mathbf{w}), \mathbf{w})$ with respect to \mathbf{w} , $b = \Delta\text{DP}$ or ΔEO
- ▶ However, optimizing $M(\hat{\theta}(\mathbf{w}), \mathbf{w})$ without any constraint on \mathbf{w} will likely result in fair but inaccurate classifier, and the optimized weights will typically not be interpretable.
- ▶ Circumventing these issues by constraining the elements of \mathbf{w} to be binary.

- ▶ Let $\bar{M}_{\mathcal{D}_{\text{val}}}^b(\hat{\theta}(\mathbf{w}), \mathbf{w})$ be a linearized approximation to $M_{\mathcal{D}_{\text{val}}}^b(\hat{\theta}(\mathbf{w}), \mathbf{w})$ about 1

Proposition

Let $\mathbf{w}_{\text{fair}} \in \{0, 1\}^N$ be a N dimensional binary vector such that its n^{th} coordinate is $w_n^{\text{fair}} = 1 - \mathbb{I}[\mathcal{I}_{b,n} > 0]$, then,

$$\mathbf{w}_{\text{fair}} = \underset{\mathbf{w} \in \{0,1\}^N}{\text{argmin}} \bar{M}_{\mathcal{D}_{\text{val}}}^b(\hat{\theta}(\mathbf{w}), \mathbf{w}) - M_{\mathcal{D}_{\text{val}}}^b(\hat{\theta}(1), 1),$$
$$\text{and } \bar{M}_{\mathcal{D}_{\text{val}}}^b(\hat{\theta}(\mathbf{w}_{\text{fair}}), \mathbf{w}_{\text{fair}}) - M_{\mathcal{D}_{\text{val}}}^b(\hat{\theta}(1), 1) \leq 0$$

- ▶ It follows that $M_{\mathcal{D}_{\text{val}}}^b(\hat{\theta}(\mathbf{w}_{\text{fair}}), \mathbf{w}_{\text{fair}}) \approx \leq M_{\mathcal{D}_{\text{val}}}^b(\hat{\theta}(1), 1)$

- ▶ Define $\mathcal{D}_- = \{\mathbf{z}_n \mid \mathbf{z}_n \in \mathcal{D} \text{ and } \mathcal{I}_{M,n} > 0\}$

$$\Rightarrow \hat{\theta}_{\text{fair}} = \hat{\theta} + \sum_{m \in \mathcal{D}_-} H^{-1} g_m$$

- ▶ In this paper, they develop an alternative iterative procedure based on the recently proposed WoodFisher approximation.
- ▶ **WoodFisher approximation**

$$\widehat{H}_{n+1}^{-1} = \widehat{H}_n^{-1} - \frac{\widehat{H}_n^{-1} \nabla_{\theta} \ell(y_{n+1}, h_{\theta}(\mathbf{x}_{n+1})) \nabla_{\theta} \ell(y_{n+1}, h_{\theta}(\mathbf{x}_{n+1}))^{\top} \widehat{H}_n^{-1}}{N + \nabla_{\theta} \ell(y_{n+1}, h_{\theta}(\mathbf{x}_{n+1}))^{\top} \widehat{H}_n^{-1} \nabla_{\theta} \ell(y_{n+1}, h_{\theta}(\mathbf{x}_{n+1}))}$$

with $\widehat{H}_0^{-1} = \lambda^{-1} I_D$, and λ a small positive scalar.

Compute Hessian-Vector Product(HVP)

► IHVP-WoodFisher approximation

⇒ purpose : $H^{-1}\mathbf{v}$

Proposition

Let $\mathbf{o}_1 = \nabla_{\theta} \ell(\mathbf{z}_1)$, $\mathbf{k}_1 = \mathbf{v}$, and N denote the number of training instances. The Hessian-vector product $H^{-1}\mathbf{v}$ is approximated by iterating through the IHVP-WoodFisher recurrence in under equation and computing \mathbf{k}_N .

$$\mathbf{o}_{n+1} = \mathbf{o}_n - \frac{\mathbf{o}_n \nabla_{\theta} \ell(\mathbf{z}_{n+1})^{\top} \mathbf{o}_n}{N + \nabla_{\theta} \ell(\mathbf{z}_{n+1})^{\top} \mathbf{o}_n}, \quad \mathbf{k}_{n+1} = \mathbf{k}_n - \frac{\mathbf{o}_n \nabla_{\theta} \ell(\mathbf{z}_{n+1})^{\top} \mathbf{k}_n}{N + \nabla_{\theta} \ell(\mathbf{z}_{n+1})^{\top} \mathbf{o}_n}$$

where, we use $\ell(\mathbf{z}_{n+1})$ as shorthand for $\ell(y_{n+1}, h_{\theta}(\mathbf{x}_{n+1}))$.

Experiment

1. Adult dataset

- Task : to predict if a person has an income above a threshold
- Sensitive attribute : **Gender** \in [Female, Male]
- Response : **Income** \in ['<=50k', '>50k']

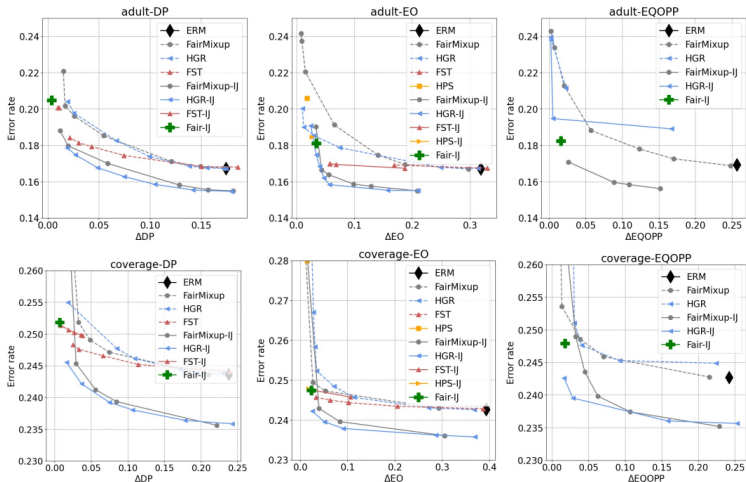
2. ACSPublicCoverage dataset

- Task : to predict whether an individual is covered by public health insurance
- Sensitive attribute : **Race** \in [white, black]
- Response : **PUBCOV**(Public health coverage) \in [0, 1]

Algorithm 1 Fair-IJ

- 1: **Input:** Pre-trained model parameters $\hat{\theta}$, training set \mathcal{D} , loss function ℓ , a validation set \mathcal{D}_{val} and a smooth surrogate to the fairness metric $b \in \{\Delta\text{DP}, \Delta\text{EO}\}$, $M_{\mathcal{D}_{\text{val}}}^b$.
 - 2: **Calculate:** $\nabla_{\theta} M(\hat{\theta}, \mathbf{1})$ using Equation 7 or Equation 8.
 - 3: **Calculate:** $\mathbf{r} = H^{-1} \nabla_{\theta} M(\hat{\theta}, \mathbf{1})$ by setting $\mathbf{k}_1 = \nabla_{\theta} M(\hat{\theta}, \mathbf{1})$ and iterating through Equation 14 for B iterations.
 - 4: **Calculate:** the fairness influence $\mathcal{I}_{b,n}$ of each training instance \mathbf{z}_n on \mathcal{D}_{val} by computing dot product between g_n and \mathbf{r} .
 - 5: **Construct:** the set \mathcal{D}_- and denote its cardinality, $|\mathcal{D}_-| = K$.
 - 6: **Initialize:** $\hat{\theta}_{\text{fair}}^0 := \hat{\theta}$
 - 7: **for** $k \in [1, \dots, K]$ **do**
 - 8: **Construct:** $\mathcal{D}_-^k = \{\mathbf{z}_n \in \mathcal{D}_- \mid \mathcal{I}_{b,n} > \mathcal{I}_{b,(K-k)}\}$, where $\mathcal{I}_{b,(K-k)}$ denotes the $(K-k)^{\text{th}}$ order statistic of the influence scores $[\mathcal{I}_{b,1}, \dots, \mathcal{I}_{b,K}]$.
 - 9: **Calculate:** $\hat{\theta}_{\text{fair}}^k$ by replacing \mathcal{D}_- with using \mathcal{D}_-^k in Equation 12.
 - 10: **If** $b_{\mathcal{D}_{\text{val}}}(\hat{\theta}_{\text{fair}}^k) < b_{\mathcal{D}_{\text{val}}}(\hat{\theta}_{\text{fair}}^{k-1})$ **set** $\hat{\theta}_{\text{fair}} := \hat{\theta}_{\text{fair}}^k$ **else set** $\hat{\theta}_{\text{fair}} := \hat{\theta}_{\text{fair}}^{k-1}$ **and break out of the for loop.**
 - 11: **end for**
 - 12: **Return:** fair model parameters $\hat{\theta}_{\text{fair}}$.
-

Result



► Accuracy and fairness Pareto frontier for the Adult and the Coverage datasets averaged over 10 runs.

► Points closer to the bottom-left achieve the best fairness/accuracy trade-off.

End