

Fool SHAP with Stealthily Biased Sampling

¹Gabriel Laberge ²Ulrich Aïvodji ³Satoshi Hara
⁴Mario Marchand ¹Foutse Khomh

¹Polytechnique Montréal

²École de Technologie Supérieure

³Osaka University

⁴Université Laval à Québec

6 avril 2023



I have a private dataset $D = \{\mathbf{x}^{(i)}\}_{i=1}^N$
and a black-box $f : \mathcal{X} \rightarrow [0, 1]$ to deploy.
The feature $x_s \in \{\text{woman}, \text{man}\}$
is sensitive.

Company

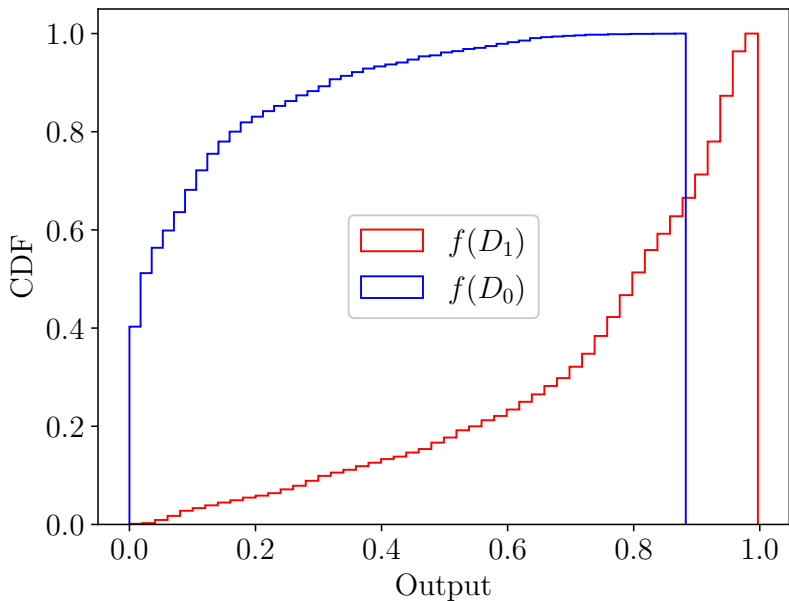
Auditor

To verify the model, we need to measure its **fairness** metrics. Can you provide access to collection of outputs $f(D_{\text{woman}}), f(D_{\text{man}})$?

$$D_{\text{woman}} = \{\mathbf{x}^{(i)} : x_s^{(i)} = \text{woman}\},$$
$$D_{\text{man}} = \{\mathbf{x}^{(i)} : x_s^{(i)} = \text{man}\}$$

Company

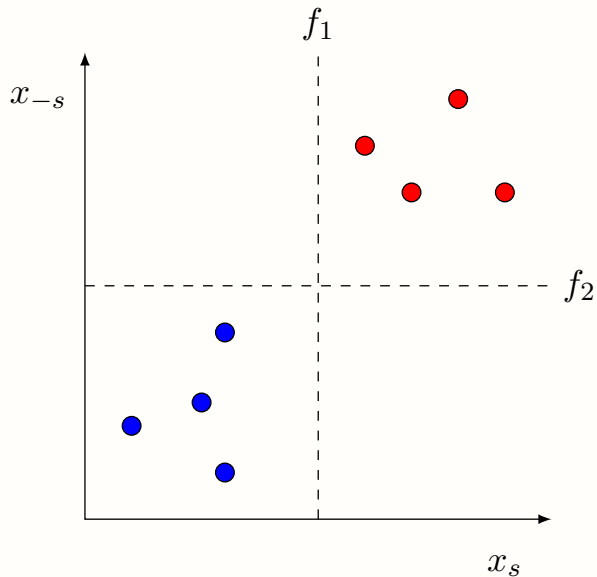
Auditor



There may be a disparity in model outcomes but that does not mean that the model is relying on the sensitive feature.
The model may rely on **meritocratic** features correlated with x_s .

Company

Auditor



To validate your argument we could compute the **Shapley Values** Φ and see which features contribute the most to the disparity.

Company

Auditor

Shapley Values to Explain Fairness

$$\sum_{i=1}^d \Phi_i(f, D_{\text{woman}}, D_{\text{man}}) = \mathbb{E}[f(\mathbf{x}) | x_s = \text{woman}] - \mathbb{E}[f(\mathbf{x}) | x_s = \text{man}]. \quad (1)$$

Shapley Values to Explain Fairness

$$\sum_{i=1}^d \Phi_i(f, D_{\text{woman}}, D_{\text{man}}) = \mathbb{E}[f(\mathbf{x}) | x_s = \text{woman}] - \mathbb{E}[f(\mathbf{x}) | x_s = \text{man}]. \quad (1)$$

Constraint

In practice, a Monte-Carlo estimate $\widehat{\Phi}(f, S_{\text{woman}}, S_{\text{man}})$ is used with two subsets $S_{\text{woman}} \subset D_{\text{woman}}$ and $S_{\text{man}} \subset D_{\text{man}}$ sampled uniformly at random.

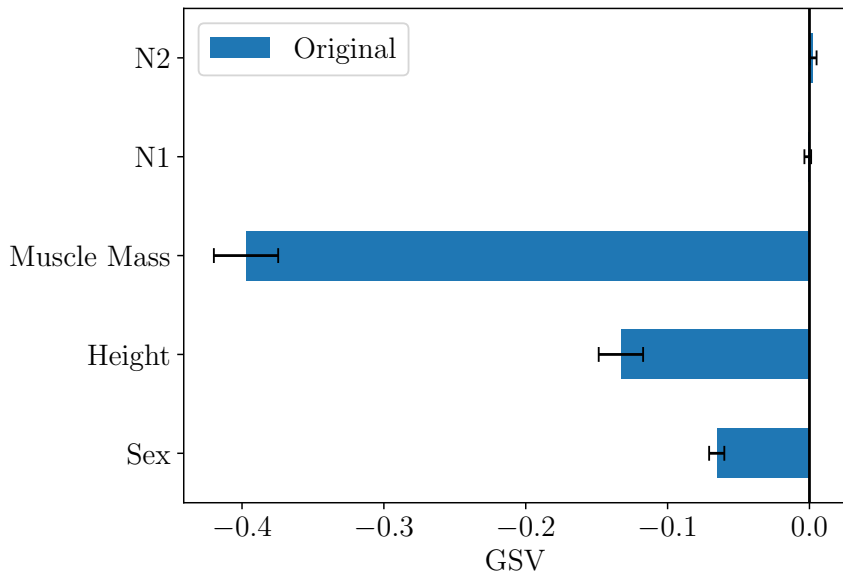
I will share with you two subsets
 $S_{\text{woman}} \subset D_{\text{woman}}$ and $S_{\text{man}} \subset D_{\text{man}}$
of size M so you can run SHAP on our
model and get $\hat{\Phi}(f, S_{\text{woman}}, S_{\text{man}})$.

Company

Auditor

Ok let's run SHAP on our own and see what we get.

Company



Ouch ! Is there a way to cherry-pick the subsets $S'_{\text{woman}}, S'_{\text{man}}$ so that $|\widehat{\Phi}_s(f, S'_{\text{woman}}, S'_{\text{man}})|$ is small and the auditor **cannot detect** the manipulation ?

Company

Detection

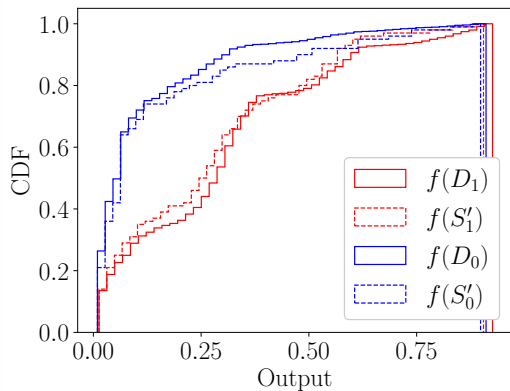
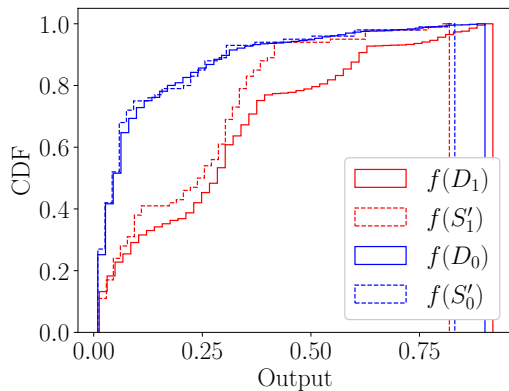
The audit already has access to $f(D_{\text{woman}})$, $f(D_{\text{man}})$. Hence they can detect the manipulation with a statistical test

$$\text{Detect_fraud}(f(D_{\text{woman}}), f(D_{\text{man}}), f(S'_{\text{woman}}), f(S'_{\text{man}}))$$

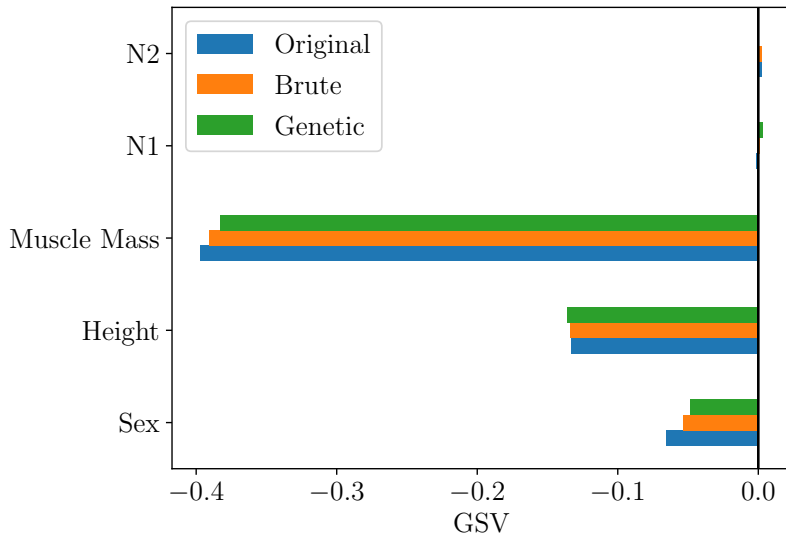
Detection

The audit already has access to $f(D_{\text{woman}}), f(D_{\text{man}})$. Hence they can detect the manipulation with a statistical test

$$\text{Detect_fraud}(f(D_{\text{woman}}), f(D_{\text{man}}), f(S'_{\text{woman}}), f(S'_{\text{man}}))$$



Baselines



Method 3 : Fool SHAP

Issues with Genetic Algorithm

- 1 Feature correlations are ignored in cross-over operation.
- 2 There is no notion of proximity to the original data.

Method 3 : Fool SHAP

Issues with Genetic Algorithm

- 1 Feature correlations are ignored in cross-over operation.
- 2 There is no notion of proximity to the original data.

Solution : Fool SHAP

- 1 Sample S'_{woman} uniformly at random.

Method 3 : Fool SHAP

Issues with Genetic Algorithm

- 1 Feature correlations are ignored in cross-over operation.
- 2 There is no notion of proximity to the original data.

Solution : Fool SHAP

- 1 Sample S'_{woman} uniformly at random.
- 2 Define $\mathcal{B} = \frac{1}{N_{\text{man}}} \sum_{\mathbf{x}^{(i)} \in D_{\text{man}}} \delta(\mathbf{x}^{(i)})$ and $\mathcal{B}' = \sum_{\mathbf{x}^{(i)} \in D_{\text{man}}} \omega_i \delta(\mathbf{x}^{(i)})$

Method 3 : Fool SHAP

Issues with Genetic Algorithm

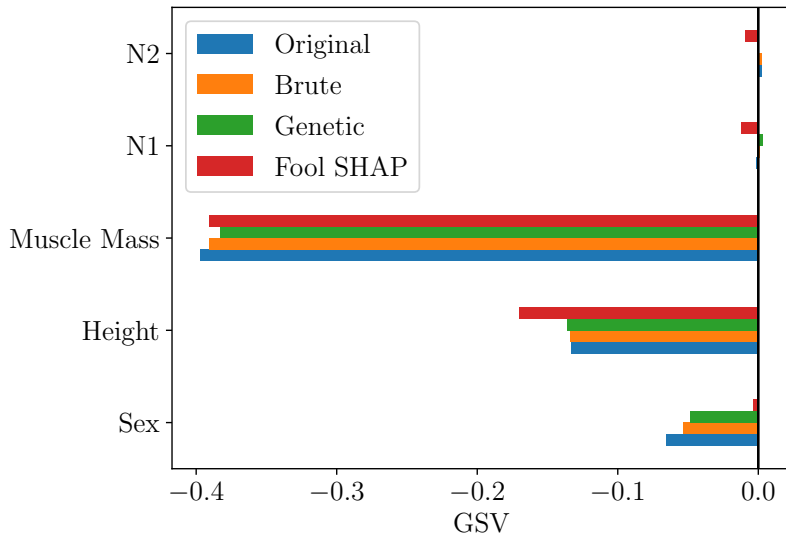
- 1 Feature correlations are ignored in cross-over operation.
- 2 There is no notion of proximity to the original data.

Solution : Fool SHAP

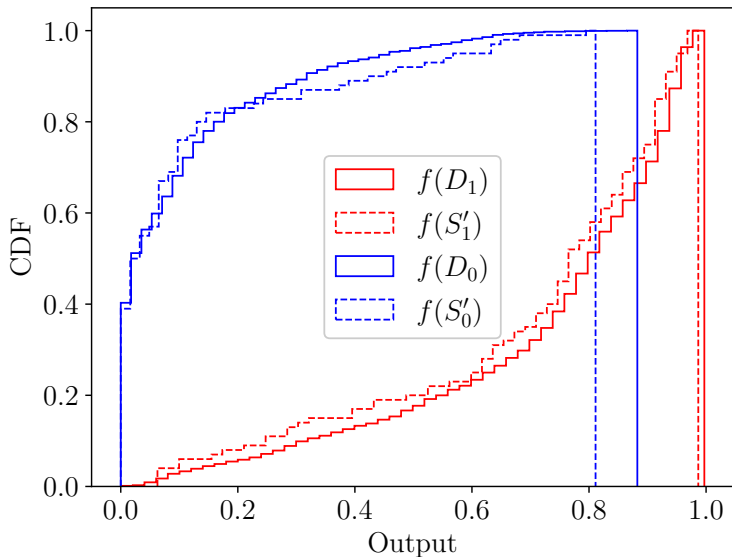
- 1 Sample S'_{woman} uniformly at random.
- 2 Define $\mathcal{B} = \frac{1}{N_{\text{man}}} \sum_{\mathbf{x}^{(i)} \in D_{\text{man}}} \delta(\mathbf{x}^{(i)})$ and $\mathcal{B}' = \sum_{\mathbf{x}^{(i)} \in D_{\text{man}}} \omega_i \delta(\mathbf{x}^{(i)})$
- 3 Optimize the weights ω such that :
 - $|\hat{\Phi}_s(f, S'_{\text{woman}}, S'_{\text{man}})|$ with $S'_{\text{man}} \sim \mathcal{B}'^M$ is small.
 - \mathcal{B}' is close to \mathcal{B} w.r.t the **Wasserstein Distance**.

Solved with a Minimum Cost Flow (MCF) Linear Program.

Method 3 Fool SHAP



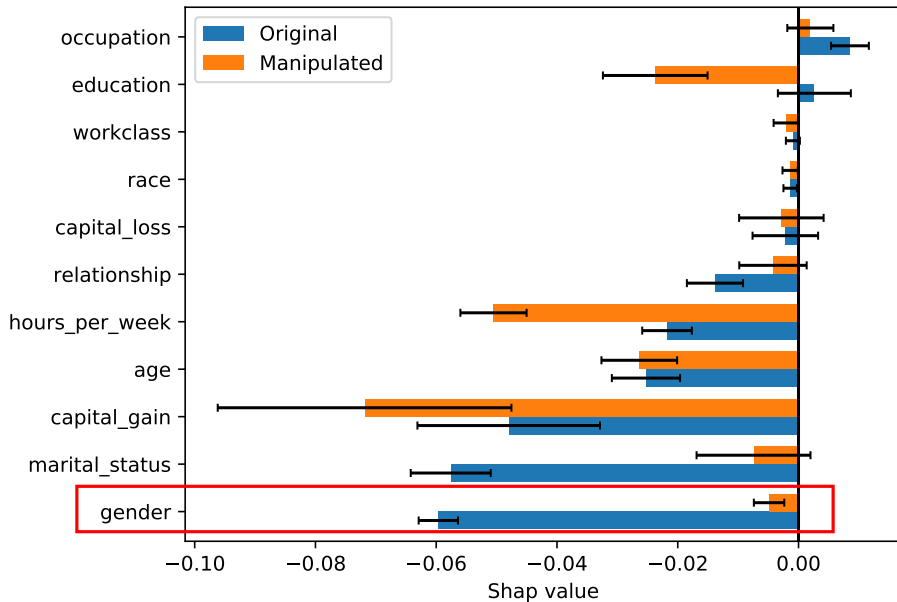
Method 3 Fool SHAP

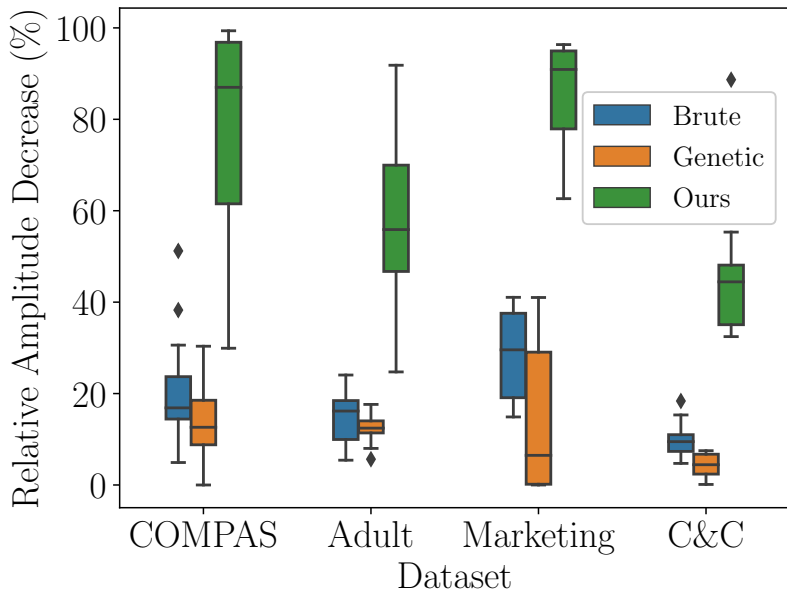


Here are the subsets S'_{woman} , S'_{man} requested.

Company

Auditor





Conclusion

Contributions

- A new and effective attack on SHAP.
- Said attacks are hard to detect by an external auditor.
- An auditor would need some access to the **input features** of the private data to circumvent the attack.

Conclusion

Contributions

- A new and effective attack on SHAP.
- Said attacks are hard to detect by an external auditor.
- An auditor would need some access to the **input features** of the private data to circumvent the attack.

Future Work

- Allow the audit to query more information about the private dataset.
- Cherry-pick S'_{woman} et S'_{man} simultaneously (Bilinear Problem).
- Apply to other measures of fairness.