# Fairwashing SHAP

## or

## Interventional and Observational Shapley Values

# Motivation

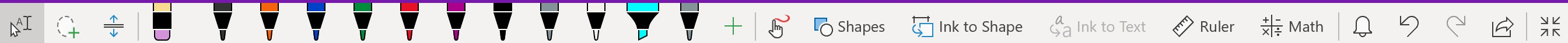Someone walks into a bank, applied for a loan, gets rejected. Why?

2 Kinds of explanations:

Recourse: What could they do to get the loan?

R\AI Developer: Why did the model say no? It is bigoted?

Focus on this

# Local Feature Importance

Given: dataset $(\underline{x}_1, y_1), \cdots, (\underline{x}_n, y_n) \overset{iid}{\sim} \mathcal{D}$ , ML model $f(\underline{x})$ , query point $\underline{q} \in \mathbb{R}^d$

*vectors are underlined*

Return: Importance of each feature $\phi_1(\underline{q}), \cdots, \phi_d(\underline{q})$

$$(\underline{x}, y_1), \cdots, (\underline{x}_n, y_n) \qquad f$$

$$\underline{q} \longrightarrow \boxed{\text{SHAP}} \longrightarrow \phi_1(\underline{q}), \cdots, \phi_d(\underline{q})$$

Shap is well used, based on nice game theory -- good about proxy vars!

# Fairwashing

Suppose $f$ is very racist.

Then $SHAP(f, q)$ should show that race is an important feature

But, we can make ML model $\tilde{f}$ such that

    ① For almost all $q$, $\tilde{f}(q) = f(q)$

    ② $SHAP(\tilde{f}, q)$ shows almost no important for race

$\tilde{f}$ is a "Fairwashed" version of $f$. It looks fair, but it ain't!

    · At least 2 papers: Umang's paper, [Slack et al. 2020]

    · But, how is this possible? What about game theory?

# Takeaway

Q: How is fairwashing possible? What about game theory?
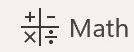
A: "Interventional" vs "Observational" Shapley Values

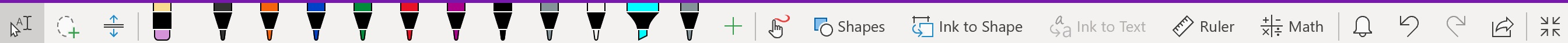What SHAP and QII do
Bad for proxys

Other papers do this.
Good for proxys.

"True to the Model or True to the Data"  (on arXiv)
Great jumping off point, easy read

# SHAPLEY VALUES

# Some notation

- We have $d$ features: $\underline{x_1}, \cdots, \underline{x_n}, \underline{q} \in \mathbb{R}^d$

- Let $S \subseteq \{1,2,3,\cdots,d\}$ be a subset of features

- Let $\bar{S}$ complement of $S$

- Let $(\underline{q})_S =$ entries of $\underline{q}$ indexed by $S$

- Notice: $\underline{q} = (\underline{q})_S + (\underline{q})_{\bar{S}}$

**Frankenstein Point**: $(\underline{q})_S + (\underline{x_j})_{\bar{S}}$

$d = 5$

$\underline{q} = [9 \quad 7 \quad 5 \quad 3 \quad 1]$

$S = \{1,2\}$

$\bar{S} = \{3,4,5\}$

$(\underline{q})_S = [9 \quad 7 \quad 0 \quad 0 \quad 0]$
$(\underline{q})_{\bar{S}} = [0 \quad 0 \quad 5 \quad 3 \quad 1]$

$\underline{x_1} = [2 \quad 2 \quad 2 \quad 2 \quad 2]$

$(\underline{q})_S + (\underline{x_1})_{\bar{S}} = [9 \quad 7 \quad 2 \quad 2 \quad 2]$

# Shapley Values

Not "Quantity of Interest"

Let $\mathcal{V}_s(q)$ be the Value of $q$ with respect to $S$.
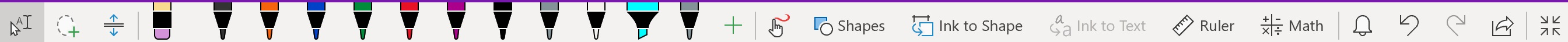The average marginal contribution of feature $i$.

$$\Phi_i(q) := \sum_{\substack{S \subseteq \{1,\dots,d\} \\ i \notin S}} \frac{1}{d} \cdot \frac{1}{\binom{d-1}{|S|}} \left( \mathcal{V}_{S \cup \{i\}}(q) - \mathcal{V}_s(q) \right)$$

Average over $S$

How much $i$ increases value, starting from $S$

① Interventional $\quad \mathcal{V}_s(q) = \mathbb{E}_{x \sim \mathcal{D}}\left[ f\left( (q)_s + (x)_{\bar{s}} \right) \right] \approx \frac{1}{n} \sum_{j=1}^{n} f\left( (q)_s + (x_j)_{\bar{s}} \right)$

- Used by SHAP (by default), QII
- Creates very fake-looking data points (ignores dependencies)

① Interventional $\quad \mathcal{V}_s(q_f) = \underset{x \sim \mathcal{D}}{\mathbb{E}} \left[ f\left( (q_f)_s + (x)_{\bar{s}} \right) \right] \approx \frac{1}{n} \sum_{j=1}^{n} f\left( (q_f)_s + (x_j)_{\bar{s}} \right)$

- Used by SHAP (by default), QII
- Creates very fake-looking data points (ignores dependencies)

② Observational $\quad \mathcal{V}_s(q_f) = \underset{x \sim \mathcal{D}}{\mathbb{E}} \left[ f\left( (q_f)_s + (x)_{\bar{s}} \right) \mid (x)_s = (q_f)_s \right]$

$$\approx \text{Average of } \left( (q_f)_s + (x_j)_{\bar{s}} \right) \text{ for all } (x_j)_s = (q_f)_s \text{ ?}$$

- Less used, but is used
- Creates real-looking data points
- Hard to compute (few points to average)

How can we compare these 2 approaches?

# Simple Linear Model

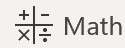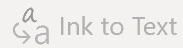Suppose $f(q_f) = \langle w, q_f \rangle + b$ for some $w \in \mathbb{R}^d$, $b \in \mathbb{R}$.

How do interventional & observational $\phi_i(q_f)$ look?

Lemma

Interventionally,

$$\phi_i(q_f) = w_i (q_i - \nu_i)$$

average of feature $i$

Takeaway: Completely ignores proxy variables

# Takeaway: Completely ignores proxy variables

$$\underline{x} = (\text{Kindness}, \text{salary}, \text{age}) \quad \text{where} \quad \text{salary} = 1000 \cdot \text{age}$$

$$\underline{w}_{BAD} = [1, 0, 1000] \quad \text{is unfair: uses age a lot}$$

$$\phi_{AGE}(\underline{q}) = 1000(q_i - \nu_i) \neq 0 \quad , \quad \text{great!}$$

$$\underline{w}_{WASHED} = [1, 1, 0] \quad \text{makes exact same predictions}$$

$$\phi_{AGE}(\underline{q}) = 0 \quad , \quad \text{terrible!}$$

**Lemma Proof:** $\phi_i(q) = w_i(q_i - \nu_i)$ interventionally

$$\mathcal{V}_S(q) = \frac{1}{n}\sum_{i=1}^{n} f\left((q)_S + (x_j)_{\bar{S}}\right) = \frac{1}{n}\sum_{i=1}^{n} \langle w, (q)_S + (x_j)_{\{i\}} + (x_j)_{\bar{S}\setminus\{i\}}\rangle + k$$

$$\mathcal{V}_{S\cup\{i\}}(q) = \frac{1}{n}\sum_{i=1}^{n} f\left((q)_{S\cup\{i\}} + (x_j)_{\bar{S}\setminus\{i\}}\right) = \frac{1}{n}\sum_{i=1}^{n} \langle w, (q)_S + (q)_{\{i\}} + (x_j)_{\bar{S}\setminus\{i\}}\rangle + k$$

$$\mathcal{V}_{S\cup\{i\}}(q) - \mathcal{V}_S(q) = \frac{1}{n}\sum_{i=1}^{n} \langle w, (q)_{\{i\}} - (x_j)_{\{i\}}\rangle$$

$$= \langle w, (q)_{\{i\}} - (\nu)_{\{i\}}\rangle$$

$$= w_i(q_i - \nu_i)$$

$$\left[0 \cdots 0 \,(q_i - \nu_i)\, 0 \cdots 0\right]$$

$$\phi_i(q) = \text{Average of } \mathcal{V}_{S\cup\{i\}}(q) - \mathcal{V}_S(q) \text{ across all} = w_i(q_i - \nu_i)$$

# Beyond Intervensional Value

Above proof does not work for observation values

$$\mathcal{V}_S(\underline{q}) = \mathop{\mathbb{E}}_{\underline{x} \sim \mathcal{D}} \left[ f((\underline{q})_S + (\underline{x})_{\bar{S}}) \mid (\underline{x})_S = (\underline{q})_S \right]$$

Only Frankenstein $\underline{q}$ with $x_j$'s that match $\underline{q}$ on $S$.
So,

$$\mathcal{V}_{S \cup \{i\}}(\underline{q}) - \mathcal{V}_S(\underline{q})$$ depends on $S$, breaking the proof.

**Claim**

If $f(\underline{q}) = \langle \underline{w}, \underline{q} \rangle + \beta$ and $x_j \overset{iid}{\sim} N(\underline{\mu}, \Sigma)$, then we can write $\phi_i(\underline{q})$ exactly (but it's ugly).

But it super depends on correlation in $\Sigma$!

# Comparison

Let

$$\underline{x}_j \sim N\left(\begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 & 1 & 0.99 \\ & & \\ & 0.99 & 1 \end{bmatrix}\right)$$

Features 2 and 3 are super correlated

$$\underline{w} = \begin{bmatrix} 1 & 2 & 3 \end{bmatrix}, \quad \underline{q} = \begin{bmatrix} 1 & 1 & 1 \end{bmatrix}$$

## Interventional

$$\phi_1(\underline{q}) = 1$$

$$\phi_2(\underline{q}) = 2 \;\Bigg]\; \text{Treated}$$

$$\phi_3(\underline{q}) = 3 \quad \begin{array}{l} \text{Very} \\ \text{Different} \end{array}$$

## Observational

$$\phi_1(\underline{q}) = 1$$

$$\phi_2(\underline{q}) \approx 2.5 \;\Bigg]\; \text{Treated}$$

$$\phi_3(\underline{q}) \approx 2.5 \quad \begin{array}{l} \text{the} \\ \text{Same!} \end{array}$$

# Interventional

$$\phi_1(q_f) = 1$$

$$\phi_2(q_f) = 2 \left.\vphantom{\begin{array}{c}a\\b\end{array}}\right] \text{Treated Very Different}$$

$$\phi_3(q_f) = 3$$

$$\underline{w} = [1 \ 0 \ 3]$$

If $w_2 = 0$,

Then $\phi_2(q_f) = 0$

Better recourse
IF

Feature are independent

# Observational

$$\phi_1(q_f) = 1$$

$$\phi_2(q_f) \approx 2.5 \left.\vphantom{\begin{array}{c}a\\b\end{array}}\right] \text{Treated the Same!}$$

$$\phi_3(q_f) \approx 2.5$$

If $w_2 = 0$

Then $\phi_2(q_f)$ may be $\neq 0$

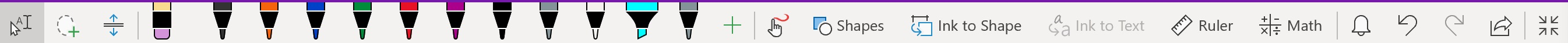If $f$ is blind to age, it possible that $\phi_{AGE}(q_f) \neq 0!$

Better recourse in my opinion
because

Features are dependant

# Conclusion

Going back to the bank example,

Both recourse and R\AI Engineering prefer observational values

But we use interventional almost always.
Also, it's harder to compute.