Carnegie Mellon University

Faith-Shap: the Faithful Shapley Interaction Index Che-Ping Tsai, Chih-Kuan Yeh, Pradeep Ravikumar Machine Learning Department, Carnegie Mellon University

TL;DR

Notations: number of input features $d \in N$, target set function $v: 2^d \to R$, maximum We propose the faithful Shapley interaction index which interaction order: $\ell \in N$, number of interactions: $d_{\ell} = \sum_{j=0}^{\ell} \binom{\ell}{j}$, the vector of interaction is the unique interaction index satisfying the original indices $\mathcal{E}(v, \ell) \in \mathbb{R}^{d_{\ell}}$, the set of subsets of size $\leq \ell: S_{\ell}$. <u>Shapley axioms while being faithful polynomial</u> approximation to the model. v_2 , and any two scalars $\alpha_1, \alpha_2 \in \mathbb{R}$, the interaction index satisfies: $\mathcal{E}(\alpha_1 v_1 + \alpha_2 v_2, \ell) = \alpha_1 \mathcal{E}(v_1, \ell) + \alpha_2 \mathcal{E}(v_2, \ell)$.

Introduction

Quantifying feature interaction is useful in some tasks.

- Question-context interaction in question answering.
- Medicine-medicine interaction in healthcare.

•Existing works extend the Shapley values to the interaction context by introducing less suitable axioms to XAI to ensure uniqueness.

- The recursive axiom in Shapley interaction indices has almost no physical meaning.
- The interaction distribution axiom in Shapley-Taylor interaction indices causes impoverished lower-order interaction attribution.

•We extend the "faithful linear approximation" property of the singleton Shapley value to "faithful polynomial approximation" for interaction indices.

- Faithfulness is an important concept in XAI.
- Together with original Shapley axioms, we prove Faith-Shap is the **unique** interaction index satisfying all these axiomatic properties.
- Similar theoretical results also apply to the Banzhaf

Preliminary: Shapley axioms

Axiom 4. (Interaction Symmetry): For any maximum interaction order $\ell \in [d]$, and for any set function $v : 2^d \mapsto \mathbb{R}$ that is symmetric with respect to elements $i, j \in [d]$, so that $v(S \cup i) = v(S \cup j)$ for any $S \subseteq [d] \setminus \{i, j\}$, the interaction index satisfies: $\mathcal{E}_{T \cup i}(v, \ell) = \mathcal{E}_{T \cup j}(v, \ell)$ for any $T \subseteq [d] \setminus \{i, j\}$ with $|T| < \ell$.

Axiom 5. (Interaction Dummy): For any maximum interaction order $\ell \in [d]$, and for any set function $v : 2^d \mapsto \mathbb{R}$ such that $v(S \cup i) = v(S)$ for some $i \in [d]$ and for all $S \subseteq [d] \setminus \{i\}$, the interaction index satisfies: $\mathcal{E}_T(v, \ell) = 0$ for all $T \in S_{\ell}$ with $i \in T$.

Axiom 6. (Interaction Efficiency): For any maximum interaction order $\ell \in [d]$, and for any set function $v : 2^d \to \mathbb{R}$, the interaction index satisfies: $\sum_{S \in S_{\ell} \setminus \emptyset} \mathcal{E}_{S}(v, \ell) = v([d]) - v(\emptyset)$ and $\mathcal{E}_{\emptyset}(v, \ell) = v(\emptyset)$.

(Faithfulness property) The singleton Shpley values are the minimizer of the following weighted linear approximation: $\min_{\mathcal{E} \in \mathbb{R}^{d+1}} \sum_{S \subseteq [d] : \, \mu(S) < \infty} \mu(S) \left(v(S) - \sum_{i \in S} \mathcal{E}_i \right)$

It has been shown that we can recover the singleton Shapley values as the solution of the weighted regression problem above by setting $\mu(S) \propto \frac{d-1}{\binom{d}{|S|} |S| (d-|S|)}$ and $\mu(\emptyset) = \mu([d]) = \infty$ [2].

Main Theoretical Results

 $\mathcal{E}(v,\ell) = \min_{v \in \mathcal{V}} \mathcal{E}(v,\ell)$

- Our theoretical results:
- (Linearity) Faith-Interaction Indices satisfy the interaction linearity axiom.
- the weighting function $\mu(S) = \mu(T)$ for all |S| = |T|.
- symmetry, and dummy axiom if and only if the weighting function satisfies:

$$\mu(S) \propto \sum_{i=|S|}^{d} \binom{d-|S|}{i-|S|} (-1)^{i-|S|} g(a,b,i), \text{ where } g(a,b,i) = \begin{cases} 1 & \text{if } i=0\\ \prod_{j=0}^{j=i-1} \frac{a(a-b)+j(b-a^2)}{a-b+j(b-a^2)} & \text{if } 1 \leq i \leq d \end{cases}$$

- $\mu([d]) = \mu(\emptyset) = \infty.$
- function is given as

$$\mu(S) \propto rac{d-1}{inom{d}{|S|} \left|S
ight| \left(d-|S|
ight)}$$

Axiom 3. (Interaction Linearity): For any maximum interaction order $\ell \in [d]$, and for any two set functions v_1 and

$$S_i$$
 s.t. $v(S) = \sum_{i \in S} \mathcal{E}_i, \ \forall S : \mu(S) = \infty.$

• We say $\mathcal{E}(v, \ell)$ are Faith-Interaction Indices if there exists a weighting function $\mu(\cdot)$ such that

$$\inf_{\mathbb{R}^{d_{\ell}}} \sum_{S \subseteq [d]} \mu(S) \left(v(S) - \sum_{T \subseteq S, T \leq \ell} \varepsilon_T \right)^2.$$

• (Symmetry) Faith-Interaction Indices satisfy the interaction symmetry axiom if and only if

• (Linearity, Symmetry, Dummy) Faith-Interaction Indices satisfy the interaction linearity,

(Efficiency) Faith-Interaction Indices satisfy the interaction efficiency axiom if and only if

• (Linearity, Symmetry, Dummy, Efficiency) Faith-Shap is the unique interaction index that satisfies the interaction linearity, symmetry, dummy, and efficiency axiom. The weighting

 $\frac{1}{|S|} \text{ for all } S \subseteq [d] \text{ with } 1 \leq |S| \leq d-1, \text{ and } \mu(\emptyset) = \mu([d]) = \infty.$

We term this unique interaction index as the Faithful Shapley Interaction index (Faith-Shap), which has the form:

$$\mathcal{E}_{S}^{F\text{-Shap}}(v,\ell) = a(v,S) + (-1)^{\ell-|S|} \frac{|S|}{\ell+|S|} \binom{\ell}{|S|} \sum_{T \supset S, |T| > \ell} \frac{\binom{|T|-1}{\ell}}{\binom{|T|+\ell-1}{\ell+|S|}} a(v,T), \ \forall S \in \mathcal{S}_{\ell}, \tag{16}$$

weighted average of discrete derivatives:

$$\mathcal{E}^{F ext{-Shap}}_{S}(v,\ell) = rac{(2k)}{(\ell)}$$

$$\binom{d}{|S_i|} |S_i| ($$

- A BERT model

Index	Sentences (bold words are the interactions with the highest (absolute) importance values)	Model Prediction	Interaction score
1	I have Never forgot this movie. All these years and it has remained in my life.	Positive	0.818
2	TWINS EFFECT is a poor film in so many respects. The only good element is that it doesn't take itself seriously	Negative	-0.375
3	I rented this movie to get an easy, entertained view of the history of Texas. I got a headache instead.	Negative	0.396
4	Truly appalling waste of space. Me and my friend tried to watch this film to its conclusion but had to switch it off about 30 minutes from the end.	Negative	0.357
5	I still remember watching Satya for the first time. I was completely blown away.	Positive	0.283

- (Example 1) While individual words "Never" and "forgot" are negative, their joint effect is positive.
- (Example 2,3,4) non-complementary interaction effect.



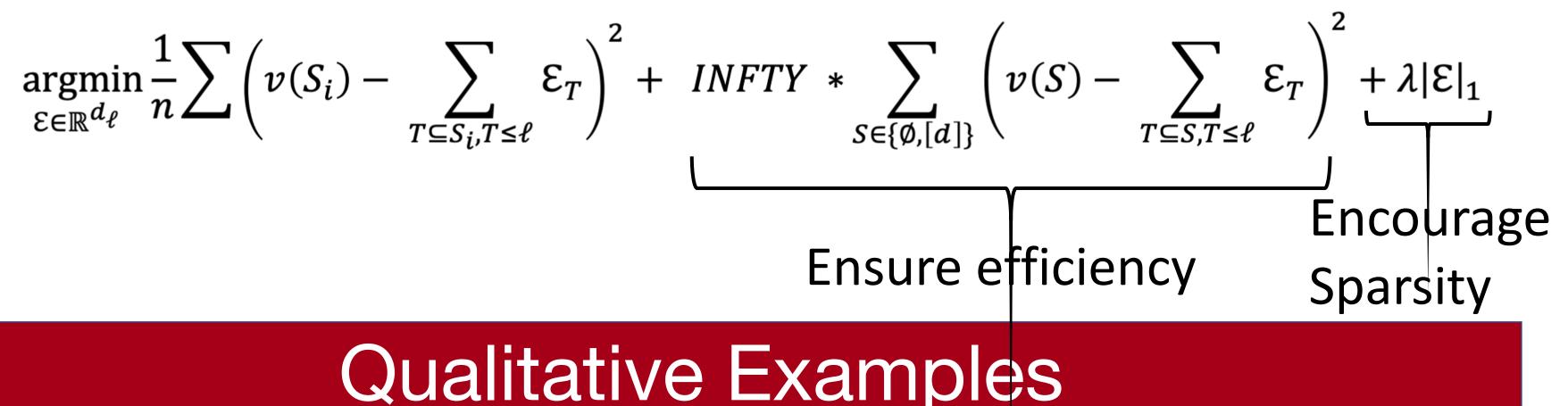
where $a(v, \cdot)$ is the Möbius transform of $v(\cdot)$. Moreover, its highest-order interaction terms can be expressed as a

 $\frac{(2\ell-1)!}{(\ell-1)!)^2} \sum_{T \in [d] \setminus S} \frac{(\ell+|T|-1)!(d-|T|-1)!}{(d+\ell-1)!} \Delta_S(v(T)) \quad \text{for all } S \in \mathcal{S}_\ell \text{ with } |S| = \ell.$ (17)

 Though the computation of the exact Faith-Shap is exponential, in practice, we use the following procedure to approximate it: e each coalition $S_i \subseteq [d]$ with probability $\mu(S_i) \propto d$

 $\frac{1}{|(d-|S_i|)}$ and compute $v(S_i)$.

2. Solve the following L1-regularized regression problem:



Explain a binary sentiment classification model on the IMDB dataset.

• Faith-Shap with the maximum interaction order $\ell = 2$.

Table 5: Top interactions of different examples on IMDB. See more results in Appendix B.

• (Example 5) Complementarity effects: words in a phrase are only meaningful when all words are present.