# On Correcting SHAP Scores

**Olivier Létoffé** [1]   **Xuanxiang Huang** [2]   **Joao Marques-Silva** [3]

## Abstract

Recent work uncovered examples of classifiers for which SHAP scores yield misleading feature attributions. While such examples might be perceived as suggesting the inadequacy of Shapley values for explainability, this paper shows that the source of the identified shortcomings of SHAP scores resides elsewhere. Concretely, the paper makes the case that the failings of SHAP scores result from the characteristic functions used in earlier works. Furthermore, the paper identifies a number of properties that characteristic functions ought to respect, and proposes several novel characteristic functions, each exhibiting one or more of the desired properties. More importantly, some of the characteristic functions proposed in this paper are guaranteed not to exhibit any of the shortcomings uncovered by earlier work. The paper also investigates the impact of the new characteristic functions on the complexity of computing SHAP scores. Finally, the paper proposes modifications to the tool SHAP to use instead one of our novel characteristic functions, thereby eliminating some of the limitations reported for SHAP scores.

## 1 Introduction

Shapley values for eXplainable AI (XAI), i.e. SHAP scores (Lundberg & Lee, 2017), are arguably among the most widely used explainability methods that target the attribution of (relative) feature importance, as exemplified by the success of the tool SHAP[1]. Despite the massive popularity of SHAP scores, some works have identified limitations with their use (Young et al., 2019; Kumar et al., 2020; Sundararajan & Najmi, 2020; Merrick & Taly, 2020; Fryer et al., 2021; Yan & Procaccia, 2021; Mothilal et al., 2021; Afchar et al., 2021; Watson et al.,

2021; Kumar et al., 2021; Campbell et al., 2022). However, most of these limitations can be attributed to the results obtained with existing tools, and not necessarily with the theoretical foundations of SHAP scores. More recent work (Huang & Marques-Silva, 2023; Huang & Marques-Silva, 2024) uncovered examples of classifiers where *exact* SHAP scores assign misleading importance to features. Namely, features having no influence in a prediction can be assigned more importance than features having the most influence in the prediction. This recent evidence should be perceived as more problematic, because it reveals limitations with the theoretical foundations of SHAP scores, and not with concrete implementations. Accordingly, these results might also be perceived as demonstrating the inadequacy of Shapley values for explainability. Nevertheless, Shapley values are of fundamental importance, not only in game theory, but also in many other domains, namely because of their intrinsic properties (Shapley, 1953).

This paper argues that the key issue with SHAP scores is not the use of Shapley values in explainability per se, and shows that the identified shortcomings of SHAP scores can be solely attributed to the characteristic functions used in earlier works (Strumbelj & Kononenko, 2010; 2014; Lundberg & Lee, 2017; Janzing et al., 2020; Sundararajan & Najmi, 2020; Arenas et al., 2021; Van den Broeck et al., 2021; 2022; Arenas et al., 2023). As noted in the recent past (Janzing et al., 2020; Sundararajan & Najmi, 2020), by changing the characteristic function, one is able to produce different sets of SHAP scores[2]. Motivated by these observations, the paper outlines fundamental properties that characteristic functions ought to exhibit in the context of XAI. Furthermore, the paper proposes several novel characteristic functions, which either respect some or all of the identified properties. In addition, the paper analyzes the impact of the novel characteristic functions on the computational complexity of computing SHAP scores, by building on recent work on the same topic (Van den Broeck et al., 2021; Arenas et al., 2021; Van den Broeck et al., 2022; Arenas et al., 2023). An indirect consequence of our work is that *corrected*

---

[1]IRIT, University of Toulouse, France; [2]CNRS@CREATE, Singapore; [3]ICREA, University of Lleida, Spain. Emails: olivier.letoffe@orange.fr, xuanxiang.huang.cs@gmail.com, jpms@icrea.cat.

[1]See https://github.com/shap/shap

[2]Unfortunately, this paper argues that past alternative proposals of characteristic functions (Janzing et al., 2020; Sundararajan & Najmi, 2020) also exhibit key limitations.

SHAP scores can be safely used for feature attribution in XAI, while offering strong guarantees regarding known shortcomings.

The paper is organized as follows. Section 2 introduces the notation and definitions used throughout the paper. Section 3 dissects some of the recently reported shortcomings with SHAP scores (Huang & Marques-Silva, 2023; Huang & Marques-Silva, 2024). The in-depth analysis of these shortcomings motivates the proposal of key properties that characteristic functions ought to exhibit. These are discussed in Section 4. Section 5 devises several novel characteristic functions, which are shown to correct some or all of the shortcomings of the characteristic functions used in earlier work. Section 6 studies the complexity of computing SHAP scores given the novel characteristic functions proposed in this paper. Section 7 outlines a simple modification to the SHAP tool (Lundberg & Lee, 2017), which corrects some of the issues of SHAP scores. Section 8 concludes the paper.

## 2 Preliminaries

**Classification problems.** Let $\mathcal{F} = \{1, \ldots, m\}$ denote a set of features and $\mathcal{K} = \{c_1, c_2, \ldots, c_K\}$ a set of classes. Each feature $i \in \mathcal{F}$ takes values from a domain $\mathbb{D}_i$. Domains can be categorical or ordinal. If ordinal, domains can be discrete or real-valued. Classes can also be categorical or ordinal. Throughout the paper domains are assumed to be discrete-valued and, unless otherwise stated, classes are assumed to be ordinal. Feature space is defined by $\mathbb{F} = \mathbb{D}_1 \times \mathbb{D}_2 \times \ldots \times \mathbb{D}_m$. The notation $\mathbf{x} = (x_1, \ldots, x_m)$ denotes an arbitrary point in feature space, where each $x_i$ is a variable taking values from $\mathbb{D}_i$. Moreover, the notation $\mathbf{v} = (v_1, \ldots, v_m)$ represents a specific point in feature space, where each $v_i$ is a constant representing one concrete value from $\mathbb{D}_i$. An *instance* denotes a pair $(\mathbf{v}, c)$, where $\mathbf{v} \in \mathbb{F}$ and $c \in \mathcal{K}$, and such that $c = \kappa(\mathbf{v})$. An ML classifier $\mathcal{M}$ is characterized by a non-constant *classification function* $\kappa$ that maps feature space $\mathbb{F}$ into the set of classes $\mathcal{K}$, i.e. $\kappa : \mathbb{F} \to \mathcal{K}$. Given the above, we associate with a classifier $\mathcal{M}$, a tuple $(\mathcal{F}, \mathbb{F}, \mathcal{K}, \kappa)$. If both $\mathbb{D}_i = \{0, 1\}, i = 1, \ldots, m$ and $\mathcal{K} = \{0, 1\}$, then the classifier is referred to as a *boolean*, in which case we use $\mathbb{B} = \{0, 1\}$. If the set of classes is ordinal but non-boolean, then the classifier is referred to as *multi-valued*. Finally, if both the domains and the set of classes are ordinal (and discrete), then the classifier is referred to as *discrete*.

**Deterministic decomposable boolean circuits (DDBCs).** For some complexity results, we will analyze DDBCs (Arenas et al., 2021; 2023)[3]. A boolean

circuit $C$ is defined on a set of (input) variables $X$ and it is represented as a directed acyclic graph, where each node is referred to as a gate, and where (i) a node with no input edges is a either a variable gate, and takes a label from $X$, or it is a constant gate, and takes a label from $\{0, 1\}$; (ii) a node with incoming edges is a either a AND, OR or NOT logic gate, where NOT gates have exactly one input; (iii) exactly one node has no output edges, and denotes the output gate of $C$. Given some circuit $C_g$, $\mathrm{var}(C_g)$ denotes the set of elements $x \in X$ such that some variable gate node of $C_g$ is labeled with $x$. A DDBC is a boolean circuit where OR gates are *deterministic* and AND gates are *decomposable*. A 2-input OR gate, $g = \mathrm{OR}(g_1, g_2)$ is deterministic if for any assignment to the inputs of the circuit, the inputs of the gate are *not* both assigned value 1. A 2-input AND gate, $g = \mathrm{AND}(g_1, g_2)$, is decomposable if $\mathrm{var}(C_{g_1})$ is disjoint from $\mathrm{var}(C_{g_2})$. It is well-known that any DDBC can be *smoothed*, i.e. all OR and AND gates can be converted to 2-input AND and OR gates (Arenas et al., 2023), in polynomial time. DDBCs generalize *deterministic decomposable negation normal form* (d-DNNF) circuits (Darwiche & Marquis, 2002). Furthermore, we consider a recent generalization of DDBCs where inputs are allowed to take multi-valued discrete values (Arenas et al., 2023).

**Selection of sets of points.** Throughout the paper, it will often be necessary to represent sets of points in feature space that are consistent with some other point in feature space with respect to the features dictated by some set of features. Accordingly, we define $\Upsilon : 2^{\mathcal{F}} \to 2^{\mathbb{F}}$ as follows[4],

$$\Upsilon(\mathcal{S}; \mathbf{v}) := \{\mathbf{x} \in \mathbb{F} \mid \wedge_{i \in \mathcal{S}} x_i = v_i\} \tag{1}$$

i.e. for some $\mathcal{S} \subseteq \mathcal{F}$, and parameterized by the point $\mathbf{v}$ in feature space, $\Upsilon(\mathcal{S}; \mathbf{v})$ denotes all the points $\mathbf{x} = (x_1, \ldots, x_m) \in \mathbb{F}$ in feature space that have in common with $\mathbf{v} = (v_1, \ldots, v_m) \in \mathbb{F}$ the values of the features specified by $\mathcal{S}$. Finally, we write $\mathbf{x}_{\mathcal{S}} = \mathbf{v}_{\mathcal{S}}$ to signify that $\mathbf{x} \in \Upsilon(\mathcal{S}; \mathbf{v})$.

**Distributions, expected value.** Throughout the paper, it is assumed a *uniform probability distribution* on each feature, and such that all features are independent. Thus, the probability of an arbitrary point in feature space becomes:

$$\mathbf{P}(\mathbf{x}) := 1/\Pi_{i \in \mathcal{F}} |\mathbb{D}_i| \tag{2}$$

That is, every point in the feature space has the same probability. The *expected value* of a classification function $\kappa$ is denoted as $\mathbf{E}[\kappa]$. Furthermore, let $\mathbf{E}[\kappa \mid \mathbf{x}_{\mathcal{S}} = \mathbf{v}_{\mathcal{S}}]$ represent the expected of $\kappa$ over points in feature space consistent with the coordinates of $\mathbf{v}$ dictated by $\mathcal{S}$, which is

---

[3]The definition of DDBC mimics the one in (Arenas et al., 2023).

[4]Parameterizations are shown as arguments after the separator ';'. However, for simplicity, we will elide parameterizations whenever these are clear from the context.

defined as follows:

$$\mathbf{E}[\kappa(\mathbf{x}) \,|\, \mathbf{x}_{\mathcal{S}} = \mathbf{v}_{\mathcal{S}}] := {}^{1}/|\Upsilon(\mathcal{S};\mathbf{v})| \sum_{\mathbf{x} \in \Upsilon(\mathcal{S};\mathbf{v})} \kappa(\mathbf{x}) \quad (3)$$

Similarly, we define,

$$\mathbf{P}(\kappa(\mathbf{x}) = c \,|\, \mathbf{x}_{\mathcal{S}} = \mathbf{v}_{\mathcal{S}}) := \qquad\qquad (4)$$
$$^{1}/|\Upsilon(\mathcal{S};\mathbf{v})| \sum_{\mathbf{x} \in \Upsilon(\mathcal{S};\mathbf{v})} \mathrm{ITE}(\kappa(\mathbf{x}) = c, 1, 0)$$

**Explanation problems.** Given a classification problem $\mathcal{M}$ and a concrete instance $(\mathbf{v}, c)$, an *explanation problem* $\mathcal{E}$ is a tuple $(\mathcal{M}, (\mathbf{v}, c))$. When describing concepts in explainability, it is to be assumed an underlying explanation problem $\mathcal{E}$, with all definitions parameterized on $\mathcal{E}$.

**Shapley values.** Shapley values were proposed in the context of game theory in the early 1950s by L. S. Shapley (Shapley, 1953). Shapley values were defined given some set $\mathcal{S}$, and a *characteristic function*, i.e. a real-valued function defined on the subsets of $\mathcal{S}$, $\upsilon : 2^{\mathcal{S}} \to \mathbb{R}$, such that $\upsilon(\emptyset) = 0$ [5]. It is well-known that Shapley values represent the *unique* function that, given $\mathcal{S}$ and $\upsilon$, respects a number of important axioms. More detail about Shapley values is available in standard references (Shapley, 1953; Dubey, 1975; Young, 1985; Roth, 1988).

**SHAP scores.** In the context of explainability, Shapley values are most often referred to as SHAP scores (Strumbelj & Kononenko, 2010; 2014; Lundberg & Lee, 2017; Arenas et al., 2021; 2023), and consider a specific characteristic function $\upsilon_e : 2^{\mathcal{F}} \to \mathbb{R}$, which is defined by,

$$\upsilon_e(\mathcal{S}; \mathcal{E}) := \mathbf{E}[\kappa(\mathbf{x}) \,|\, \mathbf{x}_{\mathcal{S}} = \mathbf{v}_{\mathcal{S}}] \qquad (5)$$

and where $\Upsilon$ (used in the definition of the expected value) is defined by (1). Thus, given a set $\mathcal{S}$ of features, $\upsilon_e(\mathcal{S}; \mathcal{E})$ represents the *e*xpected value of the classifier over the points of feature space represented by $\Upsilon(\mathcal{S}; \mathbf{v})$. The formulation presented in earlier work (Arenas et al., 2021; 2023) allows for different input distributions when computing the average values. For the purposes of this paper, it suffices to consider solely a uniform input distribution, and so the dependency on the input distribution is not accounted for. Independently of the distribution considered, it should be clear that in most cases $\upsilon_e(\emptyset) \neq 0$; this is the case for example with boolean classifiers (Arenas et al., 2021; 2023).

To simplify the notation, the following definitions are used,

$$\Delta_e(i, \mathcal{S}) := (\upsilon_e(\mathcal{S} \cup \{i\}) - \upsilon_e(\mathcal{S})) \qquad (6)$$
$$\varsigma(|\mathcal{S}|) := {}^{|\mathcal{S}|!(|\mathcal{F}| - |\mathcal{S}| - 1)!}/|\mathcal{F}|! \qquad (7)$$

Finally, let $\mathsf{Sc}_e : \mathcal{F} \to \mathbb{R}$, i.e. the SHAP score for feature $i$, be defined by,[6]

$$\mathsf{Sc}_e(i) := \sum_{\mathcal{S} \subseteq (\mathcal{F} \setminus \{i\})} \varsigma(|\mathcal{S}|) \times \Delta_e(i, \mathcal{S}) \qquad (8)$$

Given an instance $(\mathbf{v}, c)$, the SHAP score assigned to each feature measures the *contribution* of that feature with respect to the prediction. From earlier work, it is understood that a positive/negative value indicates that the feature can contribute to changing the prediction, whereas a value of 0 indicates no contribution (Strumbelj & Kononenko, 2010).

**Abductive and contrastive explanations.** Given an explanation problem $\mathcal{E}$, a *weak abductive explanation* (WAXp) is a set of features $\mathcal{S}$ such that the probability of $\kappa(\mathbf{x}) = c$ is equal to 1, when the features in $\mathcal{S}$ are assigned the values dictated by $\mathbf{v}$:

$$\mathbf{P}(\kappa(\mathbf{x}) = c \,|\, \mathbf{x}_{\mathcal{S}} = \mathbf{v}_{\mathcal{S}}) = 1 \qquad (9)$$

which implies the condition $\mathbf{E}[\kappa(\mathbf{x}) \,|\, \mathbf{x}_{\mathcal{S}} = \mathbf{v}_{\mathcal{S}}] = c$ in the case of numerical classes. An *abductive explanation* (AXp) is a subset-minimal WAXp. Similarly to the case of AXps, a *weak contrastive explanation* (WCXp) is a set of features such that the probability of $\kappa(\mathbf{x}) = c$ is less than 1, when the features not in $\mathcal{S}$ are assigned the values dictated by $\mathbf{v}$:

$$\mathbf{P}(\kappa(\mathbf{x}) = c \,|\, \mathbf{x}_{\mathcal{F} \setminus \mathcal{S}} = \mathbf{v}_{\mathcal{F} \setminus \mathcal{S}}) < 1 \qquad (10)$$

which is implied by the condition $\mathbf{E}[\kappa(\mathbf{x}) \,|\, \mathbf{x}_{\mathcal{F} \setminus \mathcal{S}} = \mathbf{v}_{\mathcal{F} \setminus \mathcal{S}}] \neq c$ in the case of numerical classes. A *contrastive explanation* (CXp) is a subset-minimal WCXp. Even though the paper defines AXps/CXps using probabilities (and expected values), these definitions are equivalent to those used in earlier works (Wäldchen et al., 2021; Marques-Silva, 2022; Darwiche, 2023). (The rationale for the alternative definitions (9) and (10) will become apparent in the following sections.)

**Feature (ir)relevancy.** The set of features that are included in at least one (abductive) explanation are defined as follows:

$$\mathfrak{F}(\mathcal{E}) := \{i \in \mathcal{X} \,|\, \mathcal{X} \in 2^{\mathcal{F}} \wedge \mathsf{AXp}(\mathcal{X})\} \qquad (11)$$

where predicate $\mathsf{AXp}(\mathcal{X})$ holds true if $\mathcal{X}$ is an AXp. (A well-known result is that $\mathfrak{F}(\mathcal{E})$ remains unchanged if CXps are used instead of AXps (Ignatiev et al., 2020), in which case predicate $\mathsf{CXp}(\mathcal{X})$ holds true if $\mathcal{X}$ is a CXp.) Finally, a feature $i \in \mathcal{F}$ is *irrelevant*, i.e. predicate $\mathsf{Irrelevant}(i)$ holds true, if $i \notin \mathfrak{F}(\mathcal{E})$; otherwise feature $i$ is *relevant*, and predicate $\mathsf{Relevant}(i)$ holds true. Clearly, given some explanation problem $\mathcal{E}$, $\forall(i \in \mathcal{F}).\mathsf{Irrelevant}(i) \leftrightarrow \neg\mathsf{Relevant}(i)$ [7].

---

[5]The original formulation also required super-additivity of the characteristic function, but that condition has been relaxed in more recent works (Dubey, 1975; Young, 1985).

[6]Throughout the paper, the definitions of $\Delta$ and $\mathsf{Sc}$ are explicitly associated with the characteristic function used in their definition.

[7]As noted earlier, the parameterization on $\mathcal{E}$ is elided.

Figure 1: Example decision tree (DT), for classifier $\kappa_{4,a}$, with target instance $((1,1,1,1),1)$.

| Classifier | $\mathsf{Sc}_e(1)$ | $\mathsf{Sc}_e(2)$ | $\mathsf{Sc}_e(3)$ | $\mathsf{Sc}_e(4)$ | Rank |
|---|---|---|---|---|---|
| $\kappa_{4,a}$ | 0.000 | 0.111 | 0.056 | -0.500 | $\langle 4,2,3,1\rangle$ |

Table 1: SHAP scores for DT in Figure 1.



| row # | $x_1$ | $x_2$ | $\kappa_1(\mathbf{x})$ |
|---|---|---|---|
| 1 | 0 | 0 | $1-6\alpha$ |
| 2 | 0 | 1 | $1+2\alpha$ |
| 3 | 1 | 0 | 1 |
| 4 | 1 | 1 | 1 |

(a) Tabular representation (TR)  (b) Decision tree (DT)

Figure 2: Simple classifier. The target instance is $((1,1),1)$.

| $\mathcal{S}$ | rows picked by $\mathcal{S}$ | $\upsilon_e(\mathcal{S})$ |
|---|---|---|
| $\emptyset$ | 1,2,3,4 | $1-\alpha$ |
| $\{1\}$ | 3,4 | 1 |
| $\{2\}$ | 2,4 | $1+\alpha$ |
| $\{1,2\}$ | 4 | 1 |

Table 2: $\upsilon_e(\mathcal{S})$ for each set $\mathcal{S}$.

## 3 Issues with SHAP Scores

Recent work (Huang & Marques-Silva, 2023; Huang & Marques-Silva, 2024) revealed a number of limitations of SHAP scores. These limitations can be categorized into two families, those occurring for boolean classifiers and those occurring for discrete (but non-boolean) classifiers.

**SHAP scores can mislead – existing example.** The example in Figure 1 is adapted from (Huang & Marques-Silva, 2024) [Fig. 8a], representing classifier $\kappa_{4,a}$. The target instance is $((1,1,1,1),1)$. It is plain that only feature 1 has influence in predicting class 1. If feature 1 takes value 1, then the prediction is *guaranteed* to be class 1. If we want to change the prediction, then we *must* change the value of feature 1. Furthermore, we need not change the value of *any* other feature, as the prediction is guaranteed to change to a value other than 1, as long as the value of feature 1 is changed. However, as can be observed in Table 1, if we compute the SHAP scores (e.g. using (8)), then the information about relative feature importance is misleading. For example, feature 4 has the largest absolute SHAP score, whereas feature 1 has a SHAP score of 0, meaning *no* importance. However, as argued above, feature 4 plays no role in setting the predicted class to 1, or in changing from that predicted class.

**SHAP scores can mislead – another example.** Simpler classifiers can be devised, which even allow selecting by how much SHAP scores mislead. For the decision tree (DT) classifier in Figure 2, Tables 2 and 3 summarize the computation of the SHAP scores, given their definition in Section 2 (using (8)). For the function not to be constant, we impose $\alpha \neq 0$. The target instance is $(\mathbf{v},c) = ((1,1),1)$. As shown below, the value of $\alpha$ defines the SHAP score of the irrelevant feature 2, whereas the Shapley value of the relevant feature 1 is always 0. Clearly, the sets of AXps/CXps are the same, i.e. $\{\{1\}\}$. Given the computed Shapley values, for $\alpha = -1/2$, we get $1 - 6\alpha = 4$ and $1 + 2\alpha = 0$, with $\mathsf{Sc}(1) = 0$ and $\mathsf{Sc}(2) = \alpha = -1/2$. However, by inspection, it is plain that feature 1 is the *only* feature that has *any* influence on the prediction, either in setting the prediction to class 1 or in changing the prediction to a class other than class 1. In contrast, feature 2 has *no* influence on the prediction, either in setting the prediction to class 1 or in changing the prediction to a class other than class 1. Thus, the computed SHAP scores would mislead a human decision maker into deeming feature 2 more important than feature 1 for the instance $((1,1,),1)$.

Earlier work (Huang & Marques-Silva, 2023; Huang & Marques-Silva, 2024) discusses several more parameterized examples of classifiers that illustrate the limitations of SHAP scores. Unfortunately, as exemplified in Appendix B, the use of baselines (Janzing et al., 2020; Sundararajan & Najmi, 2020) also reveals a number of similar limitations.

| $\mathcal{S}$ | $v_e(\mathcal{S})$ | $v_e(\mathcal{S} \cup \{1\})$ | $\Delta_e(\mathcal{S})$ | $\varsigma_e(\mathcal{S})$ | $\varsigma_e(\mathcal{S}) \times \Delta_e(\mathcal{S})$ |
|---|---|---|---|---|---|
| | | $i = 1$ | | | |
| $\emptyset$ | $1 - \alpha$ | $1$ | $\alpha$ | $1/2$ | $\alpha/2$ |
| $\{2\}$ | $1 + \alpha$ | $1$ | $-\alpha$ | $1/2$ | $-\alpha/2$ |
| | | | $\mathsf{Sc}_e(1) =$ | | $0$ |
| | | $i = 2$ | | | |
| $\mathcal{S}$ | $v_e(\mathcal{S})$ | $v_e(\mathcal{S} \cup \{2\})$ | $\Delta_e(\mathcal{S})$ | $\varsigma_e(\mathcal{S})$ | $\varsigma_e(\mathcal{S}) \times \Delta(\mathcal{S})$ |
| $\emptyset$ | $1 - \alpha$ | $1 + \alpha$ | $2\alpha$ | $1/2$ | $\alpha$ |
| $\{1\}$ | $1$ | $1$ | $0$ | $1/2$ | $0$ |
| | | | $\mathsf{Sc}_e(2) =$ | | $\alpha$ |

Table 3: Computation of Shapley values

**Discussion.** The examples above in this section reveal a number of critical issues with the existing definition of SHAP scores. Motivated by the issues with SHAP scores, several recent works (Yu et al., 2023b; Biradar et al., 2023; Yu et al., 2023a) proposed possible alternatives for feature attribution. However, the proposed alternatives are *not* themselves SHAP scores, and so do to not respect the axioms that Shapley values do.

In this paper we propose a different approach at correcting the existing issues with SHAP scores. Since the definition of SHAP scores is *unique* given some characteristic function (Shapley, 1953), we argue that the issues with SHAP scores are solely attributed to the characteristic functions used in earlier works (Strumbelj & Kononenko, 2010; 2014; Lundberg & Lee, 2017; Janzing et al., 2020; Sundararajan & Najmi, 2020; Van den Broeck et al., 2021; Arenas et al., 2021; Van den Broeck et al., 2022; Arenas et al., 2023), i.e. $v_e$ as defined by (5). The issues with $v_e$ can be categorized as follows:

1. $v_e$ is highly dependent of classes' values. In turn, this can be used to obfuscate the actual importance of features.
2. The definition of the characteristic function ignores feature (ir)relevancy, or alternatively, it ignores information about whether each set is (or is not) a (weak) AXp/CXp.
3. It is also the case that $v_e$ cannot readily be used in setting where classes are not ordinal.

(As shown in Appendix B, similar issues exist when baselines are used (Janzing et al., 2020; Sundararajan & Najmi, 2020).)

In the next section we propose properties that characteristic functions ought to exhibit, which can eliminate some or all of the issues that have been identified.

# 4 Properties of Characteristic Functions

Given the issues reported in Section 3, this section identifies properties that characteristic functions should respect. If characteristic functions fail to respect some of these properties, then the resulting SHAP scores can provide misleading information about relative feature importance.

**Weak class independence.** Let $\mathcal{M}_1 = (\mathcal{F}, \mathbb{F}, \mathcal{K}_1, \kappa_1)$ be a classifier, with domain $\mathbb{D}_i$ for each feature $i \in \mathcal{F}$. Moreover, let $\mathcal{M}_2 = (\mathcal{F}, \mathbb{F}, \mathcal{K}_2, \kappa_2)$ be another classifier, with the same domains, and with $|\mathcal{K}_1| = |\mathcal{K}_2|$. Moreover, let $\mu : \mathcal{K}_1 \to \mathcal{K}_2$ be a surjective mapping from $\mathcal{K}_1$ to $\mathcal{K}_2$, such that for any $\mathbf{x} \in \mathbb{F}$, $\kappa_2(\mathbf{x}) = \mu(\kappa_1(\mathbf{x}))$. Finally, let the target instances be $(\mathbf{v}, c)$, for $\mathcal{M}_1$, and $(\mathbf{v}, \mu(c))$ for $\mathcal{M}_2$, thus defining the explanation problems $\mathcal{E}_1 = (\mathcal{M}_1, (\mathbf{v}, c))$ and $\mathcal{E}_2 = (\mathcal{M}_2, (\mathbf{v}, \mu(c)))$. A characteristic function $v_t$ is *weakly class-independent* if, given surjective $\mu$,

$$\forall (i \in \mathcal{F}).[\mathsf{Sc}_t(i; \mathcal{E}_1) = \mathsf{Sc}_t(i; \mathcal{E}_2)]$$

**Strong class independence.** Let $\mathcal{M}_1 = (\mathcal{F}, \mathbb{F}, \mathcal{K}_1, \kappa_1)$ be a classifier, with domain $\mathbb{D}_i$ for each feature $i \in \mathcal{F}$. Moreover, let $\mathcal{M}_2 = (\mathcal{F}, \mathbb{F}, \mathcal{K}_2, \kappa_2)$ be another classifier, with the same domains. Moreover, let $\mu : \mathcal{K}_1 \to \mathcal{K}_2$ be a mapping from $\mathcal{K}_1$ to $\mathcal{K}_2$, such that for $c \in \mathcal{K}_1$, and such that,

$$\forall (b \in \mathcal{K}_1).[(b \neq c) \to (\mu(b) \neq \mu(c))]$$

Finally, let the target instances be $(\mathbf{v}, c)$, for $\mathcal{M}_1$, and $(\mathbf{v}, \mu(c))$ for $\mathcal{M}_2$, thus defining the explanation problems $\mathcal{E}_1 = (\mathcal{M}_1, (\mathbf{v}, c))$ and $\mathcal{E}_2 = (\mathcal{M}_2, (\mathbf{v}, \mu(c)))$. A characteristic function $v_t$ is *strongly class-independent* if, given $\mu$,

$$\forall (i \in \mathcal{F}).[\mathsf{Sc}_t(i; \mathcal{E}_1) = \mathsf{Sc}_t(i; \mathcal{E}_2)]$$

Given the above, the following result holds[8],

**Proposition 4.1.** *If a characteristic function is strongly class-independent, then it is weakly class-independent.*

**Compliance with feature (ir)relevancy.** Characteristic functions should respect feature (ir)relevancy, i.e. a feature is irrelevant iff its (corrected) SHAP score is 0. Formally, given a characteristic function $v_t$ is compliant with feature (ir)relevancy if,

$$\forall (i \in \mathcal{F}).\mathsf{Irrelevant}(i) \leftrightarrow (\mathsf{Sc}_t(i) = 0) \quad (12)$$

In previous work (Huang & Marques-Silva, 2023; Huang & Marques-Silva, 2024), SHAP scores are said to be *misleading* when compliance with feature (ir)relevancy is not respected. In the remainder of the paper, we assign the same meaning to the term *misleading*.

---

[8]Due to restrictions of space, most of the proofs are included in Appendix A.

**Numerical neutrality.** SHAP scores require ordinal classes, but classification often contemplates categorical classes. A characteristic function respects numerical neutrality if it can be used with both numerical and non-numerical classifiers.

**Discussion.** The properties proposed in this section target the issues reported in earlier work, where SHAP scores mislead with respect to relative feature importance. Additional properties may be devised to address any other existing issues.

## 5 New Characteristic Functions

This section proposes several new characteristic functions[9], which respect some or all of the target properties outlined in Section 4. Throughout this section, an explanation problem $\mathcal{E} = (\mathcal{M}, (\mathbf{v}, c))$ is assumed, and it is used to parameterize the proposed characteristic functions.

**Similarity function.** The new characteristic functions proposed in this paper build on a *similarity function*, $\zeta$ : $\mathbb{F} \to \{0,1\}$, that is defined as follows:

$$\zeta(\mathbf{x}; \mathcal{E}) = \begin{cases} 1 & \text{if } (\kappa(\mathbf{x}) = \kappa(\mathbf{v})) \\ 0 & \text{otherwise} \end{cases}$$

i.e. $\zeta$ takes value 1 only for the points in feature space for which the prediction matches the prediction of the target instance. Observe that $\forall(\mathcal{A} \subseteq \mathbb{F}).[\mathbf{E}[\zeta(\mathbf{x}) \,|\, \mathbf{x} \in \mathcal{A}] \in [0,1]]$. It is also plain to conclude that for $\mathcal{A}, \mathcal{B} \subseteq \mathbb{F}$, with $\mathcal{A} \subseteq \mathcal{B}$, and given $u \in \{0,1\}$, if $\mathbf{E}[\zeta(\mathbf{x} \,|\, \mathbf{x} \in \mathcal{B})] = u$ then $\mathbf{E}[\zeta(\mathbf{x} \,|\, \mathbf{x} \in \mathcal{A})] = u$, A few more properties of $\zeta$ are apparent. For $\mathcal{A} \subseteq \mathbb{F}$, $u \in \{0,1\}$, $(\mathbf{E}[\zeta(\mathbf{x}) \,|\, \mathbf{x} \in \mathcal{A}] = u) \leftrightarrow \forall(\mathbf{x} \in \mathcal{A}).[\zeta(\mathbf{x}) = u]$. As a result, it is also the case that $(\mathbf{E}[\zeta(\mathbf{x}) \,|\, \mathbf{x} \in \mathcal{A}] < 1) \leftrightarrow \exists(\mathbf{x} \in \mathcal{A}).[\zeta(\mathbf{x}) = 0]$.

**Defining the new characteristic functions.** Given the definition of the similarity function, we now introduce the following main new characteristic functions.

$$v_s(\mathcal{S}; \mathcal{E}) \;\; := \;\; \mathbf{E}[\zeta(\mathbf{x}) \,|\, \mathbf{x} \in \Upsilon(\mathcal{S}; \mathbf{v})] \tag{13}$$

$$v_a(\mathcal{S}; \mathcal{E}) \;\; := \;\; \begin{cases} 1 & \text{if } v_s(\mathcal{S}; \mathcal{E}) = 1 \\ 0 & \text{otherwise} \end{cases} \tag{14}$$

$$v_c(\mathcal{S}; \mathcal{E}) \;\; := \;\; \begin{cases} 1 & \text{if } v_s(\mathcal{F} \setminus \mathcal{S}; \mathcal{E}) < 1 \\ 0 & \text{otherwise} \end{cases} \tag{15}$$

We will refer to characteristic functions $v_e$ (see (5)), $v_s$, $v_a$, $v_c$, respectively as the *expected value*, the *similarity*, the AXp-based, and the CXp-based characteristic functions.

[9]In the rest of the paper, the symbol $v_t$ will be used to denote some concrete characteristic function distinguished by the letter $t$. The SHAP scores obtained with such characteristic function will be denoted by $\mathsf{Sc}_t$. Similarly, we will use $\Delta_t$.

Furthermore, we will introduce another characteristic function, which is shown to be tightly related with $v_a$.

$$v_n(\mathcal{S}; \mathcal{E}) \;\; := \;\; \begin{cases} 1 & \text{if } v_s(\mathcal{S}; \mathcal{E}) < 1 \\ 0 & \text{otherwise} \end{cases} \tag{16}$$

(Observe that $v_n$ can be viewed as the complement of $v_a$.)

**Basic attributes of the new characteristic functions.** We start by deriving some basic results regarding the characteristic functions $v_a$, $v_c$ and $v_n$. Throughout, it is assumed an explanation problem $\mathcal{E}$.

**Proposition 5.1.** *Given the definition of $v_a$, $v_c$ and $v_n$, then* $\mathsf{Sc}_a(i) \geq 0$, $\mathsf{Sc}_c(i) \geq 0$, *and* $\mathsf{Sc}_n(i) \leq 0$.

**Proposition 5.2.** *The following holds true:*

*1.* $\forall(\mathcal{S} \subseteq \mathcal{F}).[v_a(\mathcal{S}) = 1 \leftrightarrow \mathsf{WAXp}(\mathcal{S})]$.
*2.* $\forall(\mathcal{S} \subseteq \mathcal{F}).[v_n(\mathcal{S}) = 1 \leftrightarrow \mathsf{WCXp}(\mathcal{F} \setminus \mathcal{S})]$.
*3.* $\forall(\mathcal{S} \subseteq \mathcal{F}).[v_c(\mathcal{S}) = 1 \leftrightarrow \mathsf{WCXp}(\mathcal{S})]$.

*Proof.* We consider each case separately:

1. If $v_a(\mathcal{S}) = 1$, then, as noted earlier in the paper, $\zeta(\mathbf{x}) = 1$ for all points $\mathbf{x} \in \Upsilon(\mathcal{S})$, and so the classifier's prediction is $c$ for all points in $\Upsilon(\mathcal{S})$. Hence, by definition, $\mathcal{S}$ is a WAXp. Conversely, if $\mathcal{S}$ is an WAXp, then the prediction must be $c$ for all points $\mathbf{x}$ in $\Upsilon(\mathcal{S})$, $\forall(\mathbf{x} \in \Upsilon(\mathcal{S})).[\zeta(\mathbf{x}) = 1]$. Thus, $v_a(\mathcal{S}) = 1$.
2. If $v_n(\mathcal{S}) = 1$, then $\zeta(\mathbf{x}) \neq 1$ for some point(s) $\mathbf{x} \in \Upsilon(\mathcal{S})$, and so the classifier's prediction is not $c$ for some point(s) in $\Upsilon(\mathcal{S})$. Hence, by definition, $\mathcal{F} \setminus \mathcal{S}$ is a WCXp. Conversely, if $\mathcal{F} \setminus \mathcal{S}$ is an WCXp, then the prediction must not be $c$ for some point(s) $\mathbf{x}$ in $\Upsilon(\mathcal{S})$, i.e. $\exists(\mathbf{x} \in \Upsilon(\mathcal{S})).[\zeta(\mathbf{x} \neq 1]$. Thus, $v_n(\mathcal{S}) = 1$.
3. If $v_c(\mathcal{S}) = 1$, then $\zeta(\mathbf{x}) \neq 1$ for some point(s) $\mathbf{x} \in \Upsilon(\mathcal{F} \setminus \mathcal{S})$, and so the classifier's prediction is not $c$ for some point(s) in $\Upsilon(\mathcal{F} \setminus \mathcal{S})$. Hence, by definition, $\mathcal{S}$ is a WCXp. Conversely, if $\mathcal{S}$ is an WCXp, then the prediction must not be $c$ for some point(s) $\mathbf{x}$ in $\Upsilon(\mathcal{F} \setminus \mathcal{S})$, i.e. $\exists(\mathbf{x} \in \Upsilon(\mathcal{F} \setminus \mathcal{S})).[\zeta(\mathbf{x} \neq 1]$. Thus, $v_c(\mathcal{S}) = 1$. $\square$

**Proposition 5.3.** *The following holds true:*

*1.* $\forall(i \in \mathcal{F}).[\mathsf{Sc}_a(i) = -\mathsf{Sc}_n(i)]$.
*2.* $\forall(i \in \mathcal{F}).[\mathsf{Sc}_a(i) = \mathsf{Sc}_c(i)]$.

An immediate consequence of the results in Propositions 5.2 and 5.3, is that the complexity of computing SHAP scores $\mathsf{Sc}_t$ is the same for $t \in \{a, c, n\}$.

**Properties of the new characteristic functions.** We now prove which of the properties listed in Section 4 are respected by which characteristic functions among those proposed in this section.

It is plain that characteristic functions based on the similarity function respect numerical neutrality. Furthermore, another general result is that characteristic functions based on the similarity function guarantee strong (and so weak) class independence.

**Proposition 5.4.** *For* $t \in \{s, a, c, n\}$ *and* $i \in \mathcal{F}$, *it is the case that the characteristic function* $v_t$ *respects strong class independence.*

**Proposition 5.5.** *For* $t \in \{a, c, n\}$, *then it is the case that* $\forall(i \in \mathcal{F}).\mathsf{Irrelevant}(i) \leftrightarrow (\mathsf{Sc}_t(i) = 0)$. *Thus,* $\mathsf{Sc}_t$ *does not mislead.*

Finally, we observe that $v_s$ represents a boolean classifier, and so it exhibits the issues with SHAP scores uncovered for boolean classifiers based on $v_e$ in earlier work (Huang & Marques-Silva, 2023; Huang & Marques-Silva, 2024).

# 6 Complexity of Computing SHAP Scores

The previous section introduced novel characteristic functions that exhibit a number of desirable properties, which in turn ensure that SHAP scores will not produce misleading information. Another related question is how the novel characteristic functions impact the computional complexity of computing SHAP scores. This section starts the effort of mapping such computional complexity.

**Intractable cases.** A number of intractability results have been obtained in recent years (Van den Broeck et al., 2021; 2022). As noted earlier in the paper, for boolean functions, the similarity function does not provide any difference with respect to the original classifier. The following result is clear.

**Proposition 6.1.** *For a boolean classifier, with* $\kappa(\mathbf{v}) = 1$, *then* $\forall(\mathbf{x} \in \mathbb{F}).\zeta(\mathbf{x}; \mathcal{E}) = \kappa(\mathbf{x})$.

From Proposition 6.1 and Corollary 8 in (Van den Broeck et al., 2022), it is immediate that,

**Proposition 6.2.** *Computing SHAP scores* $\mathsf{Sc}_s$ *is #P-hard for boolean classifiers in CNF or DNF.*

Clearly, given Proposition 6.2, then the computation of SHAP scores for more complex boolean classifiers is also #P-hard.

Moreover, a key recent result regarding the computation of SHAP scores is that for the characteristic function $v_e$ there are polynomial-time algorithms for computing $\mathsf{Sc}_e$ (Arenas et al., 2021; 2023). In contrast, for characteristic functions that build on WAXps/WCXps, the computation of SHAP scores becomes NP-hard, even for d-DNNF and DDBC classifiers.

**Proposition 6.3.** *For* $t \in \{a, c, n\}$, *the computation of the*

*SHAP scores* $\mathsf{Sc}_t$ *is NP-hard for d-DNNF & DDBC classifiers.*

**Polynomial-time cases.** As shown above, the most significant tractability result that is known for $v_e$ does not hold for $v_t$, with $t \in \{a, c, n\}$. Nevertheless, some tractability results can be proved.

For classifiers represented by tabular representations (e.g. truth tables), it is simple to devise algorithms polynomial on the size of the classifier's representation (Huang & Marques-Silva, 2023).

**Proposition 6.4.** *There exist polynomial-time algorithms for computing the SHAP scores* $\mathsf{Sc}_s$, $\mathsf{Sc}_a$, $\mathsf{Sc}_c$ *for classifiers represented by tabular representations.*

Since the recent results on the tractability of computing SHAP scores for deterministic and decomposable circuits (d-DNNFs) (Arenas et al., 2021; 2023) considering boolean classifiers, then from Proposition 6.1 and (Arenas et al., 2023), it is the case that,

**Proposition 6.5.** *The computation of SHAP scores* $\mathsf{Sc}_s$ *is in P for classifiers represented by non-boolean DDBCs.*

(Observe that non-boolean d-DNNFs (Arenas et al., 2023) consider non-boolean features, but the set of classes is still binary, i.e. $\mathcal{K} = \{0, 1\}$.)

# 7 Similarity-Based SHAP

This section outlines a first step towards addressing the issues with SHAP scores reported in earlier work, and observed in the tool SHAP (Lundberg & Lee, 2017). Instead of running SHAP with the original training data and the original classifier, the similarity-based SHAP (referred to as sSHAP) replaces the original classifier by the similarity function, and reorganizes training data accordingly. In terms of running time complexity, the impact of the modifications to SHAP are negligible. More importantly, sSHAP will be approximating $\mathsf{Sc}_s$, since the underlying characteristic funtion is $v_s$. In practice, sSHAP is built on top of the SHAP tool (Lundberg & Lee, 2017).

As noted earlier in Section 5, the use of $v_s$ does not guarantee the non-existence of the issues reported in earlier work (Huang & Marques-Silva, 2024), since it is known that even boolean classifiers exhibit a number of issues related with the relative order of feature importance producing misleading information. Nevertheless, another question is whether $v_s$ can serve to correct SHAP scores (obtained with $v_e$) in classifiers for which the reported issues rely on non-boolean classification.

**Difference in SHAP scores for example DT.** Figure 3 shows the similarity function $\zeta_{4,b}$ for the classifier $\kappa_{4,a}$

Figure 3: Characteristic function $\zeta_{4,b}$ for DT of Figure 1, given instance $((1,1,1,1),1)$.

| Classifier | Sc(1) | Sc(2) | Sc(3) | Sc(4) | Rank |
|---|---|---|---|---|---|
| $\kappa_{4,a}$ | 0.000 | 0.111 | 0.056 | -0.500 | $\langle 4,2,3,1 \rangle$ |
| $\zeta_{4,b}$ | 0.500 | 0 | 0 | 0 | $\langle 1,2:3:4 \rangle$ |

Table 4: SHAP scores for $\kappa_{4,a}$ and $\zeta_{4,b}$.

shown in Figure 1. Given the obtained characteristic function, Table 4 shows the computed SHAP scores, obtained with both SHAP (Lundberg & Lee, 2017) and sSHAP. (Given the simple classifiers being considered, both SHAP and sSHAP obtain the exact SHAP scores.) For this concrete example and instance, the results confirm that the new characteristic function $\upsilon_s$ enables obtaining SHAP scores that are not misleading. As we shown next, the same situation is observed for other classifiers.

**Difference in SHAP scores for example classifiers.** To validate the improvements obtained with $\upsilon_s$ with respect to $\upsilon_e$, we studied the non-boolean classifiers reported in (Huang & Marques-Silva, 2024)[10]. For each classifier, each of the possible instances is analyzed, and the SHAP scores produced by the tools SHAP and sSHAP are recorded. If an irrelevant feature is assigned an absolute value larger than some other relevant feature, then a mismatch is declared. Table 5 summarizes the results obtained

---

[10]From (Huang & Marques-Silva, 2024), we consider (i) the two DTs of case study 2 (Fig. 3 in (Huang & Marques-Silva, 2024)), referred to as cs02a and cs02b; (ii) the two DTs of case study 3 (Fig. 5 in (Huang & Marques-Silva, 2024)), referred to as cs03a and cs03b; and (iii) the two DTs of case study 4 (Fig. 8 in (Huang & Marques-Silva, 2024)), referred to as cs04a and cs04b. Moreover, for cs02a, cs02b, cs03a and cs03b there exist 16 instances, whereas for cs04a and cs04b there exist 24 instances (because of a discrete but non-boolean domain for one of the features.)

| DT | SHAP-FRP mismatch | sSHAP-FRP mismatch |
|---|---|---|
| cs02a | 11 | 0 |
| cs02b | 4 | 0 |
| cs03a | 5 | 0 |
| cs03b | 4 | 0 |
| cs04a | 15 | 0 |
| cs04b | 4 | 0 |

Table 5: Comparison of empirical SHAP vs. empirical sSHAP.

with the two tools, where columns *SHAP-FRP mismatch* shown the number of mismatches obtained with SHAP, and column *sSHAP-FRP mismatch* shows the number of mismatches obtained with sSHAP[11]. As can be concluded, SHAP produces several mismatches. In contrast, sSHAP produces no mismatch. It should be noted that both tools are approximating the SHAP scores given the respective characteristic functions, i.e. the computed scores are not necessarily the ones dictated by (8).

As noted earlier, $\upsilon_s$ consists of replacing the original classifier by a new boolean classifier. Hence, from (Huang & Marques-Silva, 2024), such boolean classifiers can also produce misleading information. Nevertheless, given the results above and other experiments, in the cases where $\upsilon_s$ was used, we were unable to observe issue I8 (as proposed in (Huang & Marques-Silva, 2024)):

$$\forall(j \in \mathcal{F}).\,([\mathsf{Relevant}(j) \wedge (\mathsf{Sc}(j) = 0)] \vee \\ [\mathsf{Irrelevant}(j) \wedge (\mathsf{Sc}(j) \neq 0)])$$

As a result, given the observed experimental results, we make the following conjecture:

**Conjecture 1.** For the characteristic function $\upsilon_s$, issue I8 is not observed.

In the case that the above conjecture holds, another open question is whether the use of $\upsilon_s$ also prevents more flexible variants of issue I8 from being observed.

## 8   Conclusions

Recent work demonstrated the existence of classifiers for which the exact SHAP scores do not respect relative feature importance. This paper presents additional evidence, and argues that similar problems occur for SHAP scores defined in terms of baselines (Janzing et al., 2020; Sundararajan & Najmi, 2020). More importantly, the paper argues that the identified issues with SHAP scores re-

---

[11]A more extensive comparison is unrealistic at present; we would have to be able to compute exact SHAP scores, and this is only computationally feasible for very simple ML models, e.g. restricted examples of DTs.

sult from the characteristic functions used in earlier work. As a result, the paper devises several properties which characteristic functions must respect in order to compute SHAP scores that do not exhibit those issues. Complexity-wise, the paper argues that the proposed characteristic functions are as hard to compute as the characteristic functions used in earlier works studying the complexity of SHAP scores (Van den Broeck et al., 2021; Arenas et al., 2021; Van den Broeck et al., 2022; Arenas et al., 2023), or harder. Finally, the paper proposes simple modifications to the tool SHAP (Lundberg & Lee, 2017), thereby obtaining SHAP scores that respect some of the proposed properties.

## 9   Impact Statements

This paper presents work whose goal is to advance the field of Machine Learning. There are many potential societal consequences of our work, none which we feel must be specifically highlighted here.

## References

Afchar, D., Guigue, V., and Hennequin, R. Towards rigorous interpretations: a formalisation of feature attribution. In *ICML*, pp. 76–86, 2021.

Arenas, M., Barceló, P., Bertossi, L. E., and Monet, M. The tractability of SHAP-score-based explanations for classification over deterministic and decomposable boolean circuits. In *AAAI*, pp. 6670–6678, 2021.

Arenas, M., Barceló, P., Bertossi, L. E., and Monet, M. On the complexity of SHAP-score-based explanations: Tractability via knowledge compilation and non-approximability results. *J. Mach. Learn. Res.*, 24:63:1–63:58, 2023. URL http://jmlr.org/papers/v24/21-0389.html.

Biradar, G., Izza, Y., Lobo, E., Viswanathan, V., and Zick, Y. Axiomatic aggregations of abductive explanations. *CoRR*, abs/2310.03131, 2023. doi: 10.48550/arXiv.2310.03131. URL https://doi.org/10.48550/arXiv.2310.03131.

Campbell, T. W., Roder, H., Georgantas III, R. W., and Roder, J. Exact Shapley values for local and model-true explanations of decision tree ensembles. *Machine Learning with Applications*, 9:100345, 2022.

Darwiche, A. Logic for explainable AI. In *LICS*, pp. 1–11, 2023.

Darwiche, A. and Marquis, P. A knowledge compilation map. *J. Artif. Intell. Res.*, 17: 229–264, 2002. doi: 10.1613/jair.989. URL https://doi.org/10.1613/jair.989.

Dubey, P. On the uniqueness of the shapley value. *International Journal of Game Theory*, 4:131–139, 1975.

Fryer, D. V., Strümke, I., and Nguyen, H. D. Shapley values for feature selection: The good, the bad, and the axioms. *IEEE Access*, 9:144352–144360, 2021.

Huang, X. and Marques-Silva, J. The inadequacy of shapley values for explainability. *CoRR*, abs/2302.08160, 2023. doi: 10.48550/arXiv.2302.08160. URL https://doi.org/10.48550/arXiv.2302.08160.

Huang, X. and Marques-Silva, J. On the failings of shapley values for explainability. *International Journal of Approximate Reasoning*, pp. 109112, 2024. ISSN 0888-613X. doi: https://doi.org/10.1016/j.ijar.2023.109112. URL https://www.sciencedirect.com/science/article/pii

Huang, X., Cooper, M. C., Morgado, A., Planes, J., and Marques-Silva, J. Feature necessity & relevancy in ML classifier explanations. In *TACAS*, pp. 167–186, 2023.

Ignatiev, A., Narodytska, N., Asher, N., and Marques-Silva, J. From contrastive to abductive explanations and back again. In *AIxIA*, pp. 335–355, 2020.

Janzing, D., Minorics, L., and Blöbaum, P. Feature relevance quantification in explainable AI: A causal problem. In *AISTATS*, pp. 2907–2916, 2020.

Kumar, I., Scheidegger, C., Venkatasubramanian, S., and Friedler, S. A. Shapley residuals: Quantifying the limits of the Shapley value for explanations. In *NeurIPS*, pp. 26598–26608, 2021.

Kumar, I. E., Venkatasubramanian, S., Scheidegger, C., and Friedler, S. A. Problems with Shapley-value-based explanations as feature importance measures. In *ICML*, pp. 5491–5500, 2020.

Lundberg, S. M. and Lee, S. A unified approach to interpreting model predictions. In *NeurIPS*, pp. 4765–4774, 2017.

Marques-Silva, J. Logic-based explainability in machine learning. In *Reasoning Web*, pp. 24–104, 2022.

Merrick, L. and Taly, A. The explanation game: Explaining machine learning models using Shapley values. In *CDMAKE*, pp. 17–38, 2020.

Mothilal, R. K., Mahajan, D., Tan, C., and Sharma, A. Towards unifying feature attribution and counterfactual explanations: Different means to the same end. In *AIES*, pp. 652–663, 2021.

Roth, A. E. *The Shapley value: essays in honor of Lloyd S. Shapley*. Cambridge University Press, 1988.

Shapley, L. S. A value for $n$-person games. *Contributions to the Theory of Games*, 2(28):307–317, 1953.

Strumbelj, E. and Kononenko, I. An efficient explanation of individual classifications using game theory. *J. Mach. Learn. Res.*, 11:1–18, 2010. URL https://dl.acm.org/doi/10.5555/1756006.1756007.

Strumbelj, E. and Kononenko, I. Explaining prediction models and individual predictions with feature contributions. *Knowl. Inf. Syst.*, 41(3):647–665, 2014. URL https://doi.org/10.1007/s10115-013-0679-x.

Sundararajan, M. and Najmi, A. The many shapley values for model explanation. In *ICML*, pp. 9269–9278, 2020.

Valiant, L. G. The complexity of enumeration and reliability problems. *SIAM J. Comput.*, 8(3): 410–421, 1979. doi: 10.1137/0208032. URL https://doi.org/10.1137/0208032.

Van den Broeck, G., Lykov, A., Schleich, M., and Suciu, D. On the tractability of SHAP explanations. In *AAAI*, pp. 6505–6513, 2021.

Van den Broeck, G., Lykov, A., Schleich, M., and Suciu, D. On the tractability of SHAP explanations. *J. Artif. Intell. Res.*, 74:851–886, 2022. doi: 10.1613/jair.1.13283. URL https://doi.org/10.1613/jair.1.13283.

Wäldchen, S., MacDonald, J., Hauch, S., and Kutyniok, G. The computational complexity of understanding binary classifier decisions. *J. Artif. Intell. Res.*, 70: 351–387, 2021. doi: 10.1613/JAIR.1.12359. URL https://doi.org/10.1613/jair.1.12359.

Watson, D. S., Gultchin, L., Taly, A., and Floridi, L. Local explanations via necessity and sufficiency: unifying theory and practice. In *UAI*, volume 161, pp. 1382–1392, 2021.

Yan, T. and Procaccia, A. D. If you like Shapley then you'll love the core. In *AAAI*, pp. 5751–5759, 2021.

Young, H. P. Monotonic solutions of cooperative games. *International Journal of Game Theory*, 14:65–72, 1985.

Young, K., Booth, G., Simpson, B., Dutton, R., and Shrapnel, S. Deep neural network or dermatologist? *CoRR*, abs/1908.06612, 2019. URL http://arxiv.org/abs/1908.06612.

Yu, J., Farr, G., Ignatiev, A., and Stuckey, P. J. Anytime approximate formal feature attribution. *CoRR*, abs/2312.06973, 2023a. doi: 10.48550/ARXIV.2312.06973. URL https://doi.org/10.48550/arXiv.2312.06973.

Yu, J., Ignatiev, A., and Stuckey, P. J. On formal feature attribution and its approximation. *CoRR*, abs/2307.03380, 2023b. doi: 10.48550/arXiv.2307.03380. URL https://doi.org/10.48550/arXiv.2307.03380.

# A Proofs

**Proposition 4.1.** *If a characteristic function is strongly class-independent, then it is weakly class-independent.*

*Proof.* If a characteristic function is strongly class independent, it suffices to restrict the choices of $\mu$ to surjective functions to make it weakly class independent. $\square$

**Proposition 5.1.** *Given the definition of $\upsilon_a$, $\upsilon_c$ and $\upsilon_n$, then $\mathsf{Sc}_a(i) \geq 0$, $\mathsf{Sc}_c(i) \geq 0$, and $\mathsf{Sc}_n(i) \leq 0$.*

*Proof.* (Sketch) We only consider $\upsilon_a$. (The proof for $\upsilon_c$ and $\upsilon_n$ follows from Proposition 5.3.)
It is plain that $\Delta_a(i, \mathcal{S}) \in \{-1, 0, 1\}$, given the possible values that $\upsilon_a$ can take. In fact, it is the case that $\Delta_a(i, \mathcal{S}) \in \{0, 1\}$. If a set $\mathcal{S} \subseteq \mathcal{F}$ is a WAXp, then a proper superset is also a WAXp; hence it is never the case that $\Delta_a(i, \mathcal{S}) = -1$. Since every $\Delta_a(i, \mathcal{S}) \geq 0$, then $\mathsf{Sc}_a(i) \geq 0$. $\square$

**Proposition 5.3.** *The following holds true:*

1. $\forall (i \in \mathcal{F}).[\mathsf{Sc}_a(i) = -\mathsf{Sc}_n(i)]$.
2. $\forall (i \in \mathcal{F}).[\mathsf{Sc}_a(i) = \mathsf{Sc}_c(i)]$.

*Proof.* We consider each case separately:

1. By definition, it is plain that $\upsilon_a(\mathcal{S}) + \upsilon_n(\mathcal{S}) = 1$, for any $\mathcal{S} \subseteq \mathcal{F}$, because it must be the case that either $\upsilon_s(\mathcal{S}) = 1$ or $\upsilon_s(\mathcal{S}) < 1$, but not both. Given the values that $\upsilon_a(\mathcal{S})$ can take, it is also plain that $\Delta_a(i, \mathcal{S}) \in \{-1, 0, 1\}$. Moreover, if $\Delta_a(i, \mathcal{S}) = -1$, then $\Delta_n(i, \mathcal{S}) = 1$. If $\Delta_a(i, \mathcal{S}) = 1$, then $\Delta_n(i, \mathcal{S}) = -1$. Also, if $\Delta_a(i, \mathcal{S}) = 0$, then $\Delta_n(i, \mathcal{S}) = 0$. Thus, for any $i \in \mathcal{F}$ and $\mathcal{S} \subseteq \mathcal{F}$, $\Delta_n(i, \mathcal{S}) = -\Delta_a(i, \mathcal{S})$. Hence, the result follows.

2. Since $\forall (\mathcal{S} \subseteq \mathcal{F}).\mathsf{WCXp}(\mathcal{F} \setminus \mathcal{S}) \leftrightarrow \neg\mathsf{WAXp}(\mathcal{S}$, by definition, then we have $\forall (i \in \mathcal{F}), \forall (\mathcal{S} \subseteq (\mathcal{F} \setminus \{i\}))$,

$$\begin{aligned}
&\Delta_a(i, \mathcal{S}) = 1 \\
\Leftrightarrow& \neg\mathsf{WAXp}(\mathcal{S}) \wedge \mathsf{WAXp}(\mathcal{S} \cup \{i\}) \\
\Leftrightarrow& \mathsf{WCXp}(\mathcal{F} \setminus \mathcal{S}) \wedge \neg\mathsf{WAXp}(\mathcal{F} \setminus (\mathcal{S} \cup \{i\})) \\
\Leftrightarrow& \mathsf{WCXp}(\mathcal{F} \setminus \mathcal{S}) \wedge \neg\mathsf{WAXp}((\mathcal{F} \setminus \{i\}) \setminus \mathcal{S}) \\
\Leftrightarrow& \Delta_c(i, (\mathcal{F} \setminus \{i\}) \setminus \mathcal{S}) = 1
\end{aligned}$$

   Now, let $\Phi(i) := \{\mathcal{S} \subseteq (\mathcal{F} \setminus \{i\}) \,|\, \Delta_a(i, \mathcal{S}) = 1\}$. Then, by construction, $\mathsf{Sc}_a(i) = \sum_{\mathcal{S} \in \Phi(i)} \varsigma(|\mathcal{S}|)$ (because $\Delta_a = 0$ otherwise) and, by the equivalence above, $\mathsf{Sc}_c(i) = \sum_{\mathcal{S} \in \Phi(i)} \varsigma(|\mathcal{F} \setminus \{i\}) \setminus \mathcal{S}|)$. However, it is immediate to prove that $\varsigma(|\mathcal{S}|) = \varsigma(|\mathcal{F} \setminus \{i\}) \setminus \mathcal{S}|)$, and so the two sums are also equal. This proves the result. $\square$

**Proposition 5.4.** *For $t \in \{s, a, c, n\}$ and $i \in \mathcal{F}$, it is the case that the characteristic function $\upsilon_t$ respects strong class independence.*

*Proof.* For a characteristic function to respect strong class independence, the SHAP scores must not change if the classes are mapped using some function $\mu$. By hypothesis, for any point $\mathbf{x} \in \mathbb{F}$, the resulting classifier will predict $\mu(c)$ iff the original classifier predicts $c$. This means the resulting similarity functions are the same for the two classifiers, and so the SHAP scores $\mathsf{Sc}_t$, $t \in \{s, a, c, n\}$, remain unchanged. $\square$

**Proposition 5.5.** *For $t \in \{a, c, n\}$, then it is the case that $\forall (i \in \mathcal{F}).\mathsf{Irrelevant}(i) \leftrightarrow (\mathsf{Sc}_t(i) = 0)$. Thus, $\mathsf{Sc}_t$ does not mislead.*

*Proof.* First, we consider $\upsilon_a$.
Let $i \in \mathcal{F}$ be an irrelevant feature.
It is plain that $\Delta_a(i, \mathcal{S}) \in \{-1, 0, 1\}$, given the possible values that $\upsilon_a$ can take. However, as argued above, $\Delta_a(i, \mathcal{S}) \in \{0, 1\}$, since if a set $\mathcal{S} \subseteq \mathcal{F}$ is a WAXp, then a proper superset is also a WAXp; hence it is never the case that $\Delta_a(i, \mathcal{S}) = -1$. We are interested in the sets $\mathcal{S} \subseteq (\mathcal{F} \setminus \{i\})$ for which $\Delta_a(i, \mathcal{S}) = 1$, since these are the only ones that contribute to making $\mathsf{Sc}_a(i) \neq 0$. For $\Delta_a(i, \mathcal{S}) = 1$, it must be the case that $\upsilon_a(\mathcal{S}) = 0$ and $\upsilon_a(\mathcal{S} \cup \{i\}) = 1$. However, this would imply that $i$ would be included in some AXp (Huang et al., 2023). But $i$ is irrelevant, and so it is not included in any AXp. Hence, there exists no set $\mathcal{S} \subseteq (\mathcal{F} \setminus \{i\})$ such that $\Delta_a(i, \mathcal{S}) = 1$, and so $\mathsf{Sc}_a(i) = 0$.

Let $\mathsf{Sc}_a(i) = 0$.

An analysis similar to the above one allows concluding that there exist no sets $\mathcal{S}$ such that $\Delta_a(i, \mathcal{S}) = 1$. Hence, it is never the case that $v_a(\mathcal{S}) = 0$ and $v_a(\mathcal{S} \cup \{i\}) = 1$. Thus, $i$ is not included in any AXp, and so it is irrelevant.

For $v_c$ and $v_n$, it suffices to invoke Proposition 5.3; hence, the features for which $\mathsf{Sc}_a(i) = 0$ are exactly the ones for which $\mathsf{Sc}_c(i) = 0$ and $\mathsf{Sc}_n(i) = 0$.

This concludes the proof that $\mathsf{Sc}_t$, with $t \in \{a, c, n\}$, does not mislead. $\qquad\square$

**Proposition 6.2.** *Computing SHAP scores $\mathsf{Sc}_s$ is #P-hard for boolean classifiers in CNF or DNF.*

*Proof.* From (Van den Broeck et al., 2021; 2022), it is known that computing SHAP scores is polynomially equivalent to computing the expected value. In the boolean case, and so in the case of $v_s$, this is polynomially equivalent to model counting. Furthermore, model counting for DNF and CNF formulas is #P-complete (Valiant, 1979). Thus, computing the SHAP scores using $v_s$ is #P-hard. $\qquad\square$

**Proposition 6.3.** *For $t \in \{a, c, n\}$, the computation of the SHAP scores $\mathsf{Sc}_t$ is NP-hard for d-DNNF & DDBC classifiers.*

*Proof.* We reduce the problem of feature relevancy to the problem of computing the SHAP scores $\mathsf{Sc}_t$, with $t \in \{a, c, n\}$. Since feature relevancy is NP-complete for d-DNNF circuits (Huang et al., 2023), this proves that computing the SHAP scores $\mathsf{Sc}_t$, with $t \in \{a, c, n\}$ is NP-hard.

Given an explanation problem we can decide feature membership as follows. We compute the SHAP score for each feature $i \in \mathcal{F}$. Moreover, since $v_t, t \in \{a, c, n\}$ are comopliant with feature (ir)relevancy, then $\mathsf{Sc}_t(i) = 0$ iff feature $i$ is irrelevant. Hence, if we could compute the SHAP scores in polynomial-time, then we could decide feature relevancy in polynomial-time, and so computing the SHAP-scores for d-DNNFs is NP-hard.

Now, since DDBCs generalize d-DNNFs (Arenas et al., 2023), then computing the SHAP-scores for DDBCs is also NP-hard. $\qquad\square$

# B    Limitations of SHAP Scores Based on Baselines

We focus on BShap (Sundararajan & Najmi, 2020); similar analyzes could be made for other baselines (Janzing et al., 2020; Sundararajan & Najmi, 2020).

Throughout this section, the baseline is a point $\mathbf{w} \in \mathbb{F}$. Furthermore, for each $\mathcal{S} \subseteq \mathcal{F}$, let $\mathbf{x}_b^{\mathcal{S}}$ be such that $x_{b,i}^{\mathcal{S}} = \mathrm{ITE}(i \in \mathcal{S}, v_i, w_i)$.

Given $\mathbf{w} \in \mathbb{F}$, the BShap characteristic function $v_b$ is defined by $v_b(\mathcal{S}) = \kappa(\mathbf{x}_b^{\mathcal{S}})$, for $\mathcal{S} \subseteq \mathcal{F}$.

**Remarks about baselines.**    Analysis of the definition of BShap (Sundararajan & Najmi, 2020) allows proving the following results.

**Proposition B.1.** *The following holds:*

1. *BShap is only well-defined if all the domains are boolean, i.e. $\mathbb{F} = \{0, 1\}^m$.*
2. *BShap is only well-defined when $\mathbf{w} = \neg\mathbf{v}$.*

*Proof.* By contradiction, let us consider $i \in \mathcal{F}$, such that either $|\mathbb{D}_i| > 2$ or $w_i = v_i$. Then there exists a point $\mathbf{z} \in \mathbb{F}$ such that $z_i \notin \{v_i, w_i\}$. By construction, for each $\mathcal{S} \subseteq \mathcal{F}$, $\mathbf{x}_b^{\mathcal{S}}$ is different from $\mathbf{z}$. Thus, $v_b$ and so $\mathsf{Sc}_b$ do not depend on $\kappa(\mathbf{z})$. Therefore, we can use the value of $\kappa(\mathbf{z})$ to change the AXps (and CXps) without modifying the BShap scores. As there are at least $2^{(|\mathcal{F}|-1)}$ such points $\mathbf{z}$, it is plain that constructing counterexample is simple. $\qquad\square$

**BShap also misleads.**    The following notation is used $\mathcal{S} \subseteq \mathcal{F}$, let $\mathbf{v}^{\mathcal{S}}$ be defined by $\mathrm{ITE}(i \in \mathcal{S}, v_i, \neg v_i)$, with $i \in \mathcal{F}$.

For $\mathcal{S} \subseteq \mathcal{F}$, then $v_b(\mathcal{S}) = \kappa(\mathbf{v}^{\mathcal{S}})$.

**Proposition B.2.** *$v_b$ misleads.*

*Proof.* Let $\kappa(x_1, x_2) = \mathrm{ITE}(x_1 = 1, 1, 2x_2)$, and instance $(\mathbf{v}, c) = ((1, 1), 1)$.

It is plain that feature 1 influences both selecting the prediction 1 and changing the prediction to some other value. In contrast, feature 2 has not influence in either setting or changing the prediction of class 1.

It is also plain that the set of AXps is $\{\{1\}\}$, and also that $\kappa(x_1, x_2) = 1$ iff $x_1 = 1$.

However, if we compute $\mathsf{Sc}_b$, we get $\mathsf{Sc}_b(1) = 0$ and $\mathsf{Sc}_b(2) = 1$, which is of course misleading.

To confirm the SHAP scores, we proceed as follows. $\upsilon_b(\emptyset) = \kappa(0, 0) = 0$, $\upsilon_b(\{1\}) = \kappa(1, 0) = 1$, $\upsilon_b(\{2\}) = \kappa(0, 1) = 2$, and $\upsilon_b(\{1, 2\}) = \kappa(1, 1) = 1$.

Thus, $\Delta_b(1, \emptyset) = \upsilon_b(\{1\}) - \upsilon_b(\emptyset) = 1$, $\Delta_b(1, \{2\}) = \upsilon_b(\{1, 2\}) - \upsilon_b(\{2\}) = -1$, $\Delta_b(2, \emptyset) = \upsilon_b(\{2\}) - \upsilon_b(\emptyset) = 2$, $\Delta_b(2, \{1\}) = \upsilon_b(\{1, 2\}) - \upsilon_b(\{2\}) = 0$.

And finally, $\mathsf{Sc}_b(1) = (\Delta_b(1, \{2\}) + \Delta_b(1, \emptyset))/2 = 0$, $\mathsf{Sc}_b(2) = (\Delta_b(2, \{1\}) + \Delta_b(2, \emptyset))/2 = 1$. $\qquad\square$