

# The Explanation Game: Explaining Machine Learning Models Using Shapley Values

Luke Merrick<sup>1</sup> and Ankur Taly<sup>1</sup>

Fiddler Labs, Palo Alto, USA  
{luke,ankur}@fiddler.ai

**Abstract.** A number of techniques have been proposed to explain a machine learning model’s prediction by attributing it to the corresponding input features. Popular among these are techniques that apply the Shapley value method from cooperative game theory. While existing papers focus on the axiomatic motivation of Shapley values, and efficient techniques for computing them, they offer little justification for the game formulations used, and do not address the uncertainty implicit in their methods’ outputs. For instance, the popular SHAP algorithm’s formulation may give substantial attributions to features that play no role in the model. In this work, we illustrate how subtle differences in the underlying game formulations of existing methods can cause large differences in the attributions for a prediction. We then present a general game formulation that unifies existing methods, and enables straightforward confidence intervals on their attributions. Furthermore, it allows us to interpret the attributions as *contrastive explanations* of an input relative to a distribution of reference inputs. We tie this idea to classic research in cognitive psychology on contrastive explanations, and propose a conceptual framework for generating and interpreting explanations for ML models, called *formulate, approximate, explain* (FAE). We apply this framework to explain black-box models trained on two UCI datasets and a Lending Club dataset.

## 1 INTRODUCTION

Complex machine learning models are rapidly spreading to high stakes tasks such as credit scoring, underwriting, medical diagnosis, and crime prediction. Consequently, it is becoming increasingly important to interpret and explain individual model predictions to decision-makers, end-users, and regulators. A common form of model explanations are based on *feature attributions*, wherein a score (*attribution*) is ascribed to each feature in proportion to the feature’s contribution to the prediction. Over the last few years there has been a surge in feature attribution methods, with methods based on Shapley values from cooperative game theory being prominent among them [27,6,19,1,18,3,4].

Shapley values [24] provide a mathematically fair and unique method to attribute the payoff of a cooperative game to the players of the game. Recently, there have been a number of Shapley-value-based methods for attributing an ML model’s prediction to input features. Prominent among them are SHAP and KernelSHAP [19], TreeSHAP [18], QII [6], and IME [26]. In applying the Shapley values method to ML models, the key step is to setup a cooperative game whose players are the input features and whose payoff is the model prediction. Due to its strong axiomatic guarantees, the Shapley values

method is emerging as the de facto approach to feature attribution, and some researchers even speculate that it may be the only method compliant with legal regulation such as the General Data Protection Regulation’s “right to an explanation” [1].

In this work, we study several Shapley-value-based explanation techniques. Paradoxically, while all techniques lay claim to the axiomatic uniqueness of Shapley values, we discover that they yield significantly different attributions for the same input even when evaluated exactly (without approximation). In some cases, we find the attributions to be completely counter-intuitive. For instance, in Section 2, we show a simple model for which the popular SHAP method gives substantial attribution to a feature that is irrelevant to the model function. We trace this shortcoming and the differences across existing methods to the varying cooperative games formulated by the methods.<sup>1</sup> We refer to such games as *explanation games*. Unfortunately, while existing methods focus on the axiomatic motivations of Shapley values, they offer little justification for the design choices made in their explanation game formulations. The goal of this work is to shed light on these design choices, and their implications on the resulting attributions.

Our main technical result shows that various existing techniques can be unified under a common game formulation parameteric on a reference distribution. The Shapley values of this unified game formulation can be decomposed into the Shapley values of *single-reference games* that model a feature’s absence by replacing its value with the corresponding value from a specific *reference input*.

This decomposition is beneficial in two ways. First, it allows us to efficiently compute confidence intervals and other supplementary information about attributions, a notable advancement over existing methods (which lack confidence intervals even though they approximate metrics of random variables using finite samples). Second, it offers conceptual clarity. It unlocks the interpretation that attributions explain the prediction at an input *in contrast* to other reference inputs. The attributions vary across existing methods as each method chooses a different reference distribution to contrast with. We tie the idea to classic research in cognitive psychology, and propose a conceptual *formulate, approximate, explain* (FAE) framework to create Shapley-value-based *contrastive* feature attributions. The goal of the framework is to produce attributions that are not only axiomatically justified, but also relevant to the underlying explanation question.

We illustrate our ideas via case studies on models trained on two UCI datasets (Bike Sharing and Adult Income) and a Lending Club dataset. We find that in these real-world situations, explanations generated using our FAE framework uncover important patterns that previous attribution methods cannot identify. In summary, we make the following key contributions:

- We highlight several shortcomings of existing Shapley-value-based feature attribution methods (Sections 2), and analyze the root cause of these issues (Section 4.1).
- We present a novel game formulation that unifies existing methods (Section 4.2), and helps characterize their uncertainty with confidence intervals (Section 4.3).
- We offer a novel framework for creating and interpreting attributions (Section 4.4), and demonstrate its use through case studies (Section 5).

<sup>1</sup> We note that this shortcoming, and the multiplicity of game formulations has also been noted in parallel work [28,14]

$x_{male}$	$x_{lift}$	$\Pr[\mathbf{X} = \mathbf{x}]$	$f_{male}(\mathbf{x})$	$f_{both}(\mathbf{x})$
0	0	0.1	0.0	0.0
0	1	0.0	0.0	0.0
1	0	0.4	1.0	0.0
1	1	0.5	1.0	1.0

Table 1: Input distribution and model outputs for the mover hiring system example.

Payoff formulation	$f_{male}$		$f_{both}$	
	$\phi_1$ (male)	$\phi_2$ (lifting)	$\phi_1$ (male)	$\phi_2$ (lifting)
SHAP	0.05	0.05	0.028	0.472
KernelSHAP	0.10	0.00	0.050	0.450
QII	0.10	0.00	0.075	0.475
IME	0.50	0.00	0.375	0.375

Table 2: Attributions for the input  $x_{male} = 1, x_{lift} = 1$ .

## 2 A motivating example

To probe existing Shapley-value-based model explanation methods, we evaluate them on two toy models for which it is easy to intuit correct attributions. We leverage a modified version of the example provided in [6]: a system that recommends whether a moving company should hire a mover applicant. The input vector to both models comprises two binary features “is male” and “is good lifter” (denoted by  $\mathbf{x} = (x_{male}, x_{lift})$ ), and output a recommendation score between 0 (“no hire”) and 1 (“hire”). We define two models —  $f_{male}(\mathbf{x}) ::= x_{male}$  (only hire males), and  $f_{both}(\mathbf{x}) ::= x_{male} \wedge x_{lift}$  (only hire males who are good lifters). Table 1 specifies a probability distribution over the input space, along with the predictions from the two models.

Consider the input  $\mathbf{x} = (1, 1)$  (i.e. a male who is a good lifter), for which both models output a recommendation score of 1. Table 2 lists the attributions from several existing methods. Focusing on the relative attribution between  $x_{male}$  and  $x_{lift}$ , we make the following surprising observations. First, even though  $x_{lift}$  is irrelevant to  $f_{male}$ , the SHAP algorithm<sup>2</sup> results in equal attribution to both features. This contradicts our intuition around the “Dummy” axiom of Shapley values, which states that attribution to a player (feature) that never contributes to the payoff (prediction) must be zero.

Additionally, the SHAP attributions present a misleading picture from a fairness perspective:  $f_{male}$  relies solely on  $x_{male}$ , yet the attributions do not reflect this bias and instead claim that the model uses both features equally. Second, although  $f_{both}$  treats its features symmetrically, and  $\mathbf{x}$  has identical values in both its features, many of the methods considered do not provide symmetrical attributions. This again is intuitively at odds with the “Symmetry” axiom of Shapley values, which states that players (features) that always contribute equally to the payoff (prediction) must receive equal

<sup>2</sup> As defined by Equation 9 in [19].

attribution. These unintuitive behaviors surfaced by the above observations demand an in-depth study of the internal design choices of these methods. We carry out this study in Section 4.

### 3 PRELIMINARIES

#### 3.1 Additive feature attributions

*Additive feature attributions* [19] are attributions that sum to the difference between the explained model output  $f(\mathbf{x})$  and a reference output value  $\phi_0$ . In practice,  $\phi_0$  is typically an average model output or model output for a domain-specific “baseline” input (e.g. an empty string for text sentiment classification).

**Definition 1 (Additive feature attributions).** *Suppose  $f : \mathcal{X} \rightarrow \mathbb{R}$  is a model mapping an  $M$ -dimensional feature space  $\mathcal{X}$  to real-valued predictions. Additive feature attributions for  $f(\mathbf{x})$  at input  $\mathbf{x} = (x_1, \dots, x_M) \in \mathcal{X}$  comprise of a reference (or baseline) attribution  $\phi_0$  and feature attributions  $\phi = (\phi_1, \phi_2, \dots, \phi_M)$  corresponding to the  $M$  features such that  $f(\mathbf{x}) = \phi_0 + \sum_{i=1}^M \phi_i$ .*

There currently exist a number of competing methodologies for computing these attributions (see [2]). Given the difficulty of empirically evaluating attributions, several methods offer an axiomatic justification, often through the Shapley values method.

#### 3.2 Shapley values

The Shapley values method is a classic technique from game theory that fairly attributes the total payoff from a cooperative game to the game’s players [24]. Recently, this method has found numerous applications in explaining ML models (e.g. [5,19,9]).

Formally, a cooperative game is played by a set of players  $\mathcal{M} = \{1, \dots, M\}$  termed the *grand coalition*. The game is characterized by a set function  $v : 2^{\mathcal{M}} \rightarrow \mathbb{R}$  such that  $v(S)$  is the payoff for any coalition of players  $S \subseteq \mathcal{M}$ , and  $v(\emptyset) = 0$ . Shapley values are built by examining the marginal contribution of a player to an existing coalition  $S$ , i.e.,  $v(S \cup \{i\}) - v(S)$ . The Shapley value of a player  $i$ , denoted  $\phi_i(v)$ , is a certain weighted aggregation of its marginal contribution to all possible coalitions of players.

$$\phi_i(v) = \frac{1}{M} \sum_{S \subseteq \mathcal{M} \setminus \{i\}} \binom{M-1}{|S|}^{-1} (v(S \cup \{i\}) - v(S)) \quad (1)$$

The Shapley value method is the unique method satisfying four desirable axioms: *Dummy*, *Symmetry*, *Efficiency*, and *Linearity*. We informally describe the axioms in Appendix A, and refer the reader to [30] for formal definitions and proofs.

**Approximating Shapley values** Computing Shapley values involves evaluating the game payoff for every possible coalition of players. This makes the computation exponential in the number of players. For games with few players, it is possible to exactly compute the Shapley values, but for games with many players, the Shapley values can

only be approximated. Recently there has been much progress towards the efficient approximation of Shapley values. In this work we focus on a simple sampling approximation, presenting two more popular techniques in the Appendix B. We refer the reader to [20,4,1,13,3] for a fuller picture of recent advances in Shapley value approximation.

A simple sampling approximation (used by [9], among other works) relies on the fact that the Shapley value can be expressed as the expected marginal contribution a player has when players are added to a coalition in a random order. Let  $\pi(M)$  be the ordered set of permutations of  $M$ , and  $\mathcal{O}$  be an ordering randomly sampled from  $\pi(M)$ . Let  $\text{pre}_i(\mathcal{O})$  be the set of players that precede player  $i$  in  $\mathcal{O}$ . The Shapley value of player  $i$  is the expected marginal contribution of the player under all possible orderings of players.

$$\phi_i(v) = \mathbb{E}_{\mathcal{O} \sim \pi(M)} [v(\text{pre}_i(\mathcal{O}) \cup \{i\}) - v(\text{pre}_i(\mathcal{O}))] \quad (2)$$

By sampling a number of permutations and averaging the marginal contributions of each player, we can estimate this expected value for each player and approximate each player’s Shapley value.

## 4 EXPLANATION GAMES

In order to explain a model prediction with the Shapley values method, it is necessary to formulate a cooperative game with players that correspond to the features and a payoff that corresponds to the prediction. In this section, we analyze the methods examined in Section 2, and show that their surprising attributions are an artifact of their game formulations. We then discuss a unified game formulation and its decomposition to *single-reference games*, enabling conceptual clarity about the meanings of existing methods’ attributions.

**Notation** Let  $\mathcal{D}^{inp}$  be the input distribution, which characterizes the process that generates model inputs. We denote the input of an explained prediction as  $\mathbf{x} = (x_1, \dots, x_M)$  and use  $\mathbf{r}$  to denote another “reference” input. We use boldface to indicate when a variable or function is vector-valued, and capital letters for random variable inputs (although  $S$  continues to represent the set of contributing players/features). Thus,  $x_i$  is a scalar input,  $\mathbf{x}$  is an input vector, and  $\mathbf{X}$  is a *random* input vector. We use  $\mathbf{x}_S = \{x_i : i \in S\}$  to represent a sub-vector of features indexed by  $S$ . This notation is also extended to random input vectors  $\mathbf{X}$ . Lastly, we introduce the *composite input*  $\mathbf{z}(\mathbf{x}, \mathbf{r}, S)$ , which agrees with the input  $\mathbf{x}$  on all features in  $S$  and with  $\mathbf{r}$  on all features not in  $S$ . Note that  $\mathbf{z}(\mathbf{x}, \mathbf{r}, \emptyset) = \mathbf{r}$ , and  $\mathbf{z}(\mathbf{x}, \mathbf{r}, \mathcal{M}) = \mathbf{x}$ .

$$\mathbf{z}(\mathbf{x}, \mathbf{r}, S) = (z_1, z_2, \dots, z_M), \text{ where } z_i = \begin{cases} x_i & i \in S \\ r_i & i \notin S \end{cases} \quad (3)$$

### 4.1 Existing game formulations

The explanation game payoff function  $v_{\mathbf{x}}$  must be defined for every feature subset  $S$  such that  $v_{\mathbf{x}}(S)$  captures the contribution of  $\mathbf{x}_S$  to the model’s prediction. This allows

us to compute each feature’s possible marginal contributions to the prediction and derive its Shapley value (see Section 3.2).

By the definition of *additive feature attributions* (Definition 1) and the Shapley values’ Efficiency axiom, we must define  $v_{\mathbf{x}}(\mathcal{M}) ::= f(\mathbf{x}) - \phi_0$  (i.e. the payoff of the full coalition must be the difference between the explained model prediction and a baseline prediction). Although this definition is fixed, it leaves us the challenge of coming up with the payoff when some features do not contribute (that is, when they are *absent*).

We find that all existing approaches handle this feature-absent payoff by randomly sampling absent features according to a particular *reference* distribution and then computing the expected value of the prediction. The resulting game formulations differ from one another only in the reference distribution they use. Additionally, we note that in practice small samples are used to approximate the expected value present in these payoff functions. This introduces a significant source of attribution uncertainty not clearly quantified by existing work.

**Conditional distribution** The game formulation of SHAP [19], TreeSHAP [18], and [1] simulates feature absence by sampling absent features from the conditional distribution based on the values of the present (or contributing) features:

$$v_{\mathbf{x}}^{cond}(S) = \mathbb{E}_{\mathbf{R} \sim \mathcal{D}^{inp}} [f(\mathbf{z}(\mathbf{x}, \mathbf{R}, S)) \mid \mathbf{R}_S = \mathbf{x}_S] - \mathbb{E}_{\mathbf{R} \sim \mathcal{D}^{inp}} [f(\mathbf{R})] \quad (4)$$

Unfortunately, this is not a proper simulation of feature absence as it does not break correlations between features [14]. This could lead to unintuitive attributions. For instance, in the  $f_{male}$  example from Section 2, it causes the irrelevant feature  $x_{lift}$  to receive a nonzero attribution. Specifically, since the event  $x_{male} = 1$  is correlated<sup>3</sup> with  $x_{lift} = 1$ , once  $x_{lift} = 1$  is given, the expected prediction becomes 1. This causes the  $x_{lift}$  feature to have a non-zero marginal contribution (relative to when both features are absent), and therefore a nonzero Shapley value. More generally, whenever a feature is correlated with a model’s prediction on inputs drawn from  $\mathcal{D}^{inp}$ , this game formulation results in non-zero attribution to the feature regardless of whether the feature directly impacts the prediction.

**Input distribution** Another option for simulating feature absence, which is used by KernelSHAP, is to sample absent features from the corresponding marginal distribution in  $\mathcal{D}^{inp}$ :

$$v_{\mathbf{x}}^{inp}(S) = \mathbb{E}_{\mathbf{R} \sim \mathcal{D}^{inp}} [f(\mathbf{z}(\mathbf{x}, \mathbf{R}, S))] - \mathbb{E}_{\mathbf{R} \sim \mathcal{D}^{inp}} [f(\mathbf{R})] \quad (5)$$

Since this formulation breaks correlation with the contributing features, it ensures irrelevant features receive no attribution (e.g. no attribution to  $x_{lift}$  when explaining  $f_{male}(1, 1) = 1$ ). We formally describe this property via the *Insensitivity* axiom in Section 4.2.

<sup>3</sup> In this context, *correlation* refers to general statistical dependence, not just a nonzero Pearson correlation coefficient.

We note that this formulation is still subject to artifacts of the input distribution, as evident from the asymmetrical attributions when explaining the prediction  $f_{both}(1, 1) = 1$  (see Table 2). The features receive different attributions because they have different marginal distributions in  $\mathcal{D}^{inp}$ , not because they impact the model differently.

**Joint-marginal distribution** QII [6] simulates feature absence by sampling absent features one at a time from their own univariate marginal distributions. In addition to breaking correlation with the contributing features, this breaks correlation between absent features as well. Formally, the QII formulation uses a distribution we term the “joint-marginal” distribution ( $\mathcal{D}^{J.M.}$ ), where:

$$\Pr_{X \sim \mathcal{D}^{J.M.}} [X = (x_1, \dots, x_M)] = \prod_{i=1}^M \Pr_{X_i \sim \mathcal{D}^{inp}} [X_i = x_i]$$

The joint-marginal formulation  $v_{\mathbf{x}}^{J.M.}$  is similar to  $v_{\mathbf{x}}^{inp}$ , except that the reference distribution is  $\mathcal{D}^{J.M.}$  instead of  $\mathcal{D}^{inp}$ :

$$v_{\mathbf{x}}^{J.M.}(S) = \mathbb{E}_{\mathbf{R} \sim \mathcal{D}^{J.M.}} [f(\mathbf{z}(\mathbf{x}, \mathbf{R}, S))] - \mathbb{E}_{\mathbf{R} \sim \mathcal{D}^{J.M.}} [f(\mathbf{R})] \quad (6)$$

Unfortunately, like  $v_{\mathbf{x}}^{inp}$ , this game formulation is also tied to the input distribution and under-attributes features that take on common values in the background data. This is evident from the attributions for the  $f_{both}$  model shown in Table 2.

**Uniform distribution** The last formulation we study from the prior art simulates feature absence by drawing values from a uniform distribution  $\mathcal{U}$  over the entire input space, as in IME [26].<sup>4</sup> Completely ignoring the input distribution, this payoff  $v_{\mathbf{x}}^{unif}$  considers all possible feature values (edge-cases and common cases) with equal weighting.

$$v_{\mathbf{x}}^{unif}(S) = \mathbb{E}_{\mathbf{R} \sim \mathcal{U}} [f(\mathbf{z}(\mathbf{x}, \mathbf{R}, S))] - \mathbb{E}_{\mathbf{R} \sim \mathcal{U}} [f(\mathbf{R})] \quad (7)$$

In Table 2, we see that this formulation yields intuitively correct attributions for  $f_{male}$  and  $f_{both}$ . However, the uniform distribution can sample so heavily from irrelevant outlier regions of  $\mathcal{X}$  that relevant patterns of model behavior become masked (we study the importance of *relevant references* both theoretically in Section 4.4 and empirically in Section 5).

## 4.2 A unified formulation

We observe that the existing game formulations  $v_{\mathbf{x}}^{inp}$ ,  $v_{\mathbf{x}}^{J.M.}$ , and  $v_{\mathbf{x}}^{unif}$  can be unified as a single game formulation  $v_{\mathbf{x}, \mathcal{D}^{ref}}$  that is parameterized by a reference distribution  $\mathcal{D}^{ref}$ .

$$v_{\mathbf{x}, \mathcal{D}^{ref}}(S) = \mathbb{E}_{\mathbf{R} \sim \mathcal{D}^{ref}} [f(\mathbf{z}(\mathbf{x}, \mathbf{R}, S))] - \mathbb{E}_{\mathbf{R} \sim \mathcal{D}^{ref}} [f(\mathbf{R})] \quad (8)$$

<sup>4</sup> It is somewhat unclear whether IME proposes  $\mathcal{U}$  or  $\mathcal{D}^{inp}$ , as [26] assumes  $\mathcal{D}^{inp} = \mathcal{U}$ , while [27] calls for values to be sampled from  $\mathcal{X}$  “at random.”

For instance, the formulation for KernelSHAP is recovered when  $\mathcal{D}^{ref} = \mathcal{D}^{inp}$ , and QII is recovered when  $\mathcal{D}^{ref} = \mathcal{D}^{J.M.}$ . In the rest of this section, we discuss several properties of this general formulation that help us better understand its attributions. Notably, the formulation  $v_{\mathbf{x}}^{cond}$  cannot be expressed in this framework; we discuss the reason for this later in this section.

**A decomposition in terms of single-reference games** We now introduce *single-reference games*, a conceptual building block that helps us interpret the Shapley values of the  $v_{\mathbf{x}, \mathcal{D}^{ref}}$  game. A single-reference game  $v_{\mathbf{x}, \mathbf{r}}$  simulates feature absence by replacing the feature value with the value from a specific reference input  $\mathbf{r}$ :

$$v_{\mathbf{x}, \mathbf{r}}(S) = f(z(\mathbf{x}, \mathbf{r}, S)) - f(\mathbf{r}) \quad (9)$$

The attributions from a single-reference game explain the difference between the prediction for the input and the prediction for the reference (i.e.  $\sum_i \phi_i(v_{\mathbf{x}, \mathbf{r}}) = v_{\mathbf{x}, \mathbf{r}}(\mathcal{M}) = f(\mathbf{x}) - f(\mathbf{r})$ , and  $\phi_0 = f(\mathbf{r})$ ). Computing attributions relative to a single reference point (also referred to as a “baseline”) is common to several other methods [29, 25, 7, 3]. However, while those works seek a neutral “informationless” reference (e.g. an all-black image for image models), we find it beneficial to consider arbitrary references and interpret the resulting attributions relative to the reference. We develop this idea further in our FAE framework (see Section 4.4).

We now state Proposition 1, which shows how the Shapley values of  $v_{\mathbf{x}, \mathcal{D}^{ref}}$  can be expressed as the expected Shapley values of a (randomized) single-reference game  $v_{\mathbf{x}, \mathbf{R}}$ , where  $\mathbf{R} \sim \mathcal{D}$ . The proof (provided in Appendix C) follows from the Shapley values’ Linearity axiom and the linearity of expectation.

**Proposition 1.**  $\phi(v_{\mathbf{x}, \mathcal{D}^{ref}}) = \mathbb{E}_{\mathbf{R} \sim \mathcal{D}^{ref}} [\phi(v_{\mathbf{x}, \mathbf{R}})]$

Proposition 1 brings conceptual clarity and practical improvements (confidence intervals and supplementary metrics) to existing methods. It shows that the attributions from existing games ( $v_{\mathbf{x}}^{inp}$ ,  $v_{\mathbf{x}}^{J.M.}$ , and  $v_{\mathbf{x}}^{unif}$ ) are in fact differently weighted aggregations of attributions from a space of single-reference games. For instance,  $v_{\mathbf{x}}^{unif}$  weighs attributions relative to all reference points equally, while  $v_{\mathbf{x}}^{inp}$  weighs them using the input distribution  $\mathcal{D}^{inp}$ .

**Insensitivity axiom** We show that attributions from the game  $v_{\mathbf{x}, \mathcal{D}^{ref}}$  satisfy the *Insensitivity* axiom from [29], which states that a feature that is mathematically irrelevant to the model must receive zero attribution. Formally, a feature  $i$  is irrelevant to a model  $f$  if for any input, changing the feature does not change the model output. That is,  $\forall \mathbf{x}, \mathbf{r} \in \mathcal{X} : \mathbf{x}_{\mathcal{M} \setminus \{i\}} = \mathbf{r}_{\mathcal{M} \setminus \{i\}} \implies f(\mathbf{x}) = f(\mathbf{r})$ .

**Proposition 2.** *If a feature  $i$  is irrelevant to a model  $f$  then  $\phi_i(v_{\mathbf{x}, \mathcal{D}^{ref}}) = 0$  for all distributions  $\mathcal{D}^{ref}$ .*

Notably, the  $v_{\mathbf{x}}^{cond}$  formulation does not obey the Insensitivity axiom (a counter-example being the  $f_{male}$  attributions from Section 2). Accordingly, our general formulation (Equation 7) cannot express this formulation. In the rest of the paper, we focus on game formulations that satisfy the Insensitivity axiom. We refer to [28] for a comprehensive analysis of the axiomatic guarantees of various game formulations.

### 4.3 Confidence intervals on attributions

Existing game formulations involve computing an expected value (over a reference distribution) in every invocation of the payoff function. In practice, this expectation is approximated via sampling, which introduces uncertainty. The original formulations of these games do not lend themselves well to quantify such uncertainty. We show that by leveraging our unified game formulation, one can efficiently quantify the uncertainty using confidence intervals (CIs).

Our decomposition in Proposition 1 shows that the attributions themselves can be expressed as an expectation over (deterministic) Shapley value attributions from a distribution of single-reference games. Consequently, we can quantify attribution uncertainty by estimating the standard error of the mean (SEM) across a sample of Shapley values from single-reference games. In terms of the sample standard deviation (SSD), 95% CIs on the mean attribution ( $\bar{\phi}$ ) from a sample of size  $N$  are given by

$$\bar{\phi} \pm \frac{1.96 \times \text{SSD}(\{\phi(v_{\mathbf{x}, \mathbf{r}_i})\}_{i=1}^N)}{\sqrt{N}} \quad (10)$$

We note that while one could use bootstrap [8] to obtain CIs, the SEM approach is more efficient as it requires no additional Shapley value computations.

**A unified CI** As discussed in Section 3.2, often the large number of features (players) in an explanation game necessitates the approximation of Shapley values. The approximation may involve random sampling, which incurs its own uncertainty. In what follows, we derive a general SEM-based CI that quantifies the combined uncertainty from sampling-based approximations of Shapley values and the sampling of references.

Let us consider a generic estimator  $\hat{\phi}_i^{(\mathbf{G})}(v_{\mathbf{x}, \mathbf{r}})$  parameterized by some random sample  $\mathbf{G}$ . An example of such an approach is the feature ordering based approximation of Equation 2, for which  $\mathbf{G} = (\mathbf{O}_j)_{j=1}^k$  represents a random sample of feature orderings, and:

$$\hat{\phi}_i^{(\mathbf{G})}(v_{\mathbf{x}, \mathbf{r}}) = \frac{1}{k} \sum_{j=1}^k v(\text{pre}_i(\mathbf{O}_j) \cup \{i\}) - v(\text{pre}_i(\mathbf{O}_j))$$

As long as the generic  $\hat{\phi}_i^{(\mathbf{G})}$  is an unbiased estimator (like the feature ordering estimator of Equation 2), and  $\mathbf{G}$  and  $\mathbf{R} \sim \mathcal{D}^{ref}$  are sampled independently from one another, we can derive a unified CI using the SEM. By the estimator's unbiasedness and Proposition 1, the Shapley value attributions can be expressed as:

$$\phi_i(v_{\mathbf{x}, \mathcal{D}^{ref}}) = \mathbb{E}_{\mathbf{R}} \mathbb{E}_{\mathbf{G}} \left[ \hat{\phi}_i^{(\mathbf{G})}(v_{\mathbf{x}, \mathbf{R}}) \right] \quad (11)$$

Since  $\mathbf{G}$  is independent of  $\mathbf{R}$ , this expectation can be Monte Carlo estimated using the sample mean of the sequence  $\left( \hat{\phi}_i^{(\mathbf{g}_j)}(v_{\mathbf{x}, \mathbf{r}_j}) \right)_{j=1}^N$  (where  $(\mathbf{g}_j, \mathbf{r}_j)_{j=1}^N$  is a joint sample of  $(\mathbf{G}, \mathbf{R})$ ). As the attribution recovered by this estimation is simply the mean of a sample from a random variable, its uncertainty can be quantified by estimating the SEM. In

terms of the sample standard deviation, 95% CIs on the mean attribution ( $\bar{\phi}$ ) from a sample of size  $N$  are given by:

$$\bar{\phi} \pm \frac{1.96 \times \text{SSD} \left( \left( \hat{\phi}_i^{(g_j)}(v_{\mathbf{x}, \mathbf{r}_j}) \right)_{j=1}^N \right)}{\sqrt{N}} \quad (12)$$

#### 4.4 Formulate, Approximate, Explain

So far we studied the explanation game formulations used by existing methods, and noted how the formulations impact the resulting Shapley value attributions. We show that the attributions explain a prediction *in contrast* to a distribution of references; see Proposition 1. Existing methods differ in the attribution they produce because each of them picks a different reference distribution to contrast with. We also proposed a mechanism to quantify the approximation uncertainty incurred in computing attributions.

We now put these ideas together in a single conceptual framework *formulate, approximate, explain* (FAE). Our key insight is that rather than viewing the reference distribution as an implementation detail of the explanation method, it must be made a first-class argument to the framework. That is, the references must be consciously chosen by the explainee to obtain a specific *contrastive explanation*.

Our emphasis on treating attributions as contrastive explanations stems from cognitive psychology. Several works in cognitive psychology argue that humans frame explanations of surprising outcomes by contrasting them with to one or more normal outcomes [15,21,22,10,11,17,12]. In our setting, the normal outcomes are the reference predictions that the input prediction is contrasted with. The attributions essentially explain what drives the prediction at hand away from the reference predictions. The choice of references may depend on the context of the question, and may vary across explainers and explainees [15]. Moreover, it is important for the references to be *relevant* to the input at hand [11]. For instance, if we are explaining why an auto-grading software assigns a B+ to a student’s submission, it would be proper to contrast with the submissions that were graded as A- (next higher grade after B+), instead of contrasting with the entire pool of submissions.

**Formulate** The mandate of the Formulate step is to *generate a contrastive question that specifies one or more relevant references*. The question pins down the distribution  $\mathcal{D}^{ref}$  of the chosen references. For instance, in the grading example above, the references would be all submissions obtaining an A- grade.

**Approximate** Once a meaningful contrastive question and its corresponding reference distribution  $\mathcal{D}^{ref}$  has been formulated, we consider the distribution of single-reference games whose references are drawn from  $\mathcal{D}^{ref}$ , and approximate the Shapley values of these games. Formally, we approximate the distribution of the random-valued attribution vector  $\Phi_{\mathbf{x}, \mathbf{R}} = \phi(v_{\mathbf{x}, \mathbf{R}})$ , where  $\mathbf{R} \sim \mathcal{D}^{ref}$ . This involves two steps: (1) sampling a sequence of references  $(\mathbf{r}_i)_{i=1}^N$  from  $\mathbf{R} \sim \mathcal{D}^{ref}$ , and (2) approximating the Shapley value of the single-reference games relative each to reference in  $(\mathbf{r}_i)_{i=1}^N$ . This yields a

sequence of approximated Shapley values. It is important to account for the uncertainty resulting from sampling in steps (1) and (2), and quantify it in the Explain step.

**Explain** In the final step, we must summarize the sampled Shapley value vectors (drawn from  $\Phi_{\mathbf{x}, \mathbf{R}}$ ) obtained from the Approximate step. One simple summarization would be the presentation of a few representative samples, in the style of the SP-LIME algorithm [23]. Another simple summarization is the sample mean, which approximates  $\mathbb{E}[\Phi_{\mathbf{x}, \mathbf{R}}]$ , and is equivalent to the attributions from the unified explanation game  $v_{\mathbf{x}, \mathcal{D}^{ref}}$ . This is the summarization used by existing Shapley-value-based explanation methods. When using the sample mean, the framework of Section 4.3 can be used to quantify the uncertainty from sampling. In addition, one must be careful that the mean does not hide important information. For instance, a feature’s attributions may have opposite signs relative to different references. Averaging these attributions will cause them to cancel each other out, yielding a small mean that incorrectly suggests that the feature is unimportant. We discuss a concrete example of this in Section 5.1. At the very least, we recommend confirming through visualization and summary statistics like variance and interquartile range that the mean is a good summarization, before relying upon it. We discuss a clustering based summarization method in Section 5 while leaving further research on faithful summarization methods to future work.

## 5 CASE STUDIES

In this section we apply the FAE framework to LightGBM [16] Gradient Boosted Decision Trees (GBDT) models trained on real data: the UCI Bike Sharing and Adult Income datasets, and a Lending Club dataset.<sup>5</sup> For parsimony, we analyze models that use only five features; complete model details are provided in Appendix D. For the Bike Sharing model, we explain a randomly selected prediction of 210 rentals for a certain hour. For the Adult Income model, we explain a counter-intuitively low prediction for an individual with high *education-num*. For the Lending Club model, we explain a counter-intuitive rejection (assuming a threshold that accepts 15% of loan applications) for a high-income borrower. In the rest of this section, we present a selection of the results across all three models, while the full set of results are provided in Appendix E.

### 5.1 Shortcomings of existing methods

Recall from Section 4.2 that the attributions from existing methods amount to computing the mean attribution for a distribution of single-reference games  $v_{\mathbf{x}, \mathbf{R}}$ , where the reference  $\mathbf{R}$  is sampled from a certain distribution. The choice of distribution varies across the methods, which in turns leads to very different attribution. This is illustrated in Table 3 for the Bike sharing model.

<sup>5</sup> In Bike Sharing we model hourly bike rentals from temporal and weather features, in Adult Income we model whether an adult earns more than \$50,000 annually, and in Lending Club we model whether a borrower will default on a loan.

	Game	Avg. Prediction ( $\phi_0$ )	hr	temp	work.	hum	season
$v_{\mathbf{x}}^{inp}$	151		3	47	1	7	2
$v_{\mathbf{x}}^{J.M.}$	141		6	50	1	9	3
$v_{\mathbf{x}}^{unif}$	128		3	60	3	12	3

Table 3: Bike Sharing comparison of mean attributions. 95% CIs ranged from  $\pm 0.4$  (*hum* in  $\mathcal{D}^{inp}$  and  $\mathcal{D}^{J.M.}$ ) to  $\pm 2.5$  (*hr* in  $\mathcal{D}^{inp}$  and  $\mathcal{D}^{J.M.}$ ).

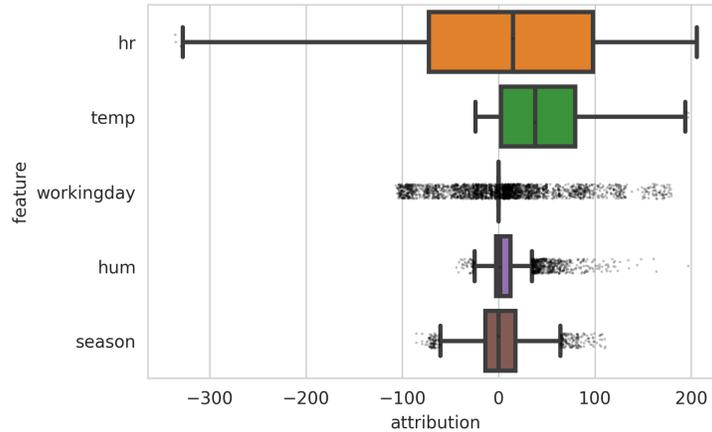


Fig. 1: Distribution of single-reference game attributions relative to the data distribution ( $\mathcal{D}^{inp}$ ) for the Bike Sharing example.

**Misleading means** In Section 4.4, we discussed that the mean attribution can potentially be a misleading summarization. Here, we illustrate this using the attributions from the KernelSHAP game  $v_{\mathbf{x}}^{inp}$  for the Bike Sharing example; see Table 3. The mean attribution to the feature *hr* is tiny, suggesting that the feature has little impact. However, the distribution of single-reference game attributions (Figure 1) reveals a large spread centered close to zero. In fact, we find that by absolute value *hr* receives the largest attribution in over 60% of the single-reference games. Consequently, only examining the mean of the distribution may be misleading.

**Unquantified uncertainty** Lack of uncertainty quantification in existing techniques can result in misleading attributions. For instance, taking the mean attribution of 100 randomly-sampled Bike Sharing single-reference games<sup>6</sup> gives *hr* an attribution of -10 and *workingday* an attribution of 8. Without any sense of uncertainty, we do not know how accurate these (estimated) attributions are. The 95% confidence intervals

<sup>6</sup> The official implementation of KernelSHAP [19] raises a warning if over 100 references are used.

Game	Size	Avg. Prediction ( $\phi_0$ )	rel.	cap.	edu.	mar.	age
$v_{\mathbf{x}}^{inp}$	-	0.24	-0.04	-0.03	-0.01	-0.10	-0.00
$v_{\mathbf{x}}^{J.M.}$	-	0.19	-0.02	-0.03	-0.01	-0.08	0.01
$v_{\mathbf{x}}^{unif}$	-	0.82	0.01	-0.79	0.02	-0.03	0.04
Cluster 1	10.2%	0.67	-0.15	-0.01	-0.15	-0.28	-0.02
Cluster 2	55.3%	0.04	0.01	0.00	0.00	-0.01	0.02
Cluster 3	4.4%	0.99	-0.04	-0.70	-0.06	-0.12	-0.01
Cluster 4	28.0%	0.31	-0.09	0.00	0.08	-0.21	-0.03
Cluster 5	2.1%	0.67	-0.04	0.01	-0.47	-0.14	0.03

Table 4: Adult Income comparison of mean attributions from existing game formulations (top) and from clusters obtained from k-means clustering of the single-reference game attributions relative to the input distribution (bottom). 95% CIs ranged from  $\pm 0.0004$  (Cluster 2, *relationship*) to  $\pm 0.0115$  (Cluster 5, *marital-status* and *age*).

Game	Avg. Prediction ( $\phi_0$ )	fico.	addr.	inc.	acc.	dti
$v_{\mathbf{x}}^{inp}$	0.14	0.00	0.03	0.00	0.10	0.00
$v_{\mathbf{x}, \mathcal{D}^{accept}}$	0.05	0.02	0.04	0.02	0.11	0.03

Table 5: Lending Club comparison of mean attributions from the game  $v_{\mathbf{x}}$  (relative to the data distribution  $\mathcal{D}^{inp}$ ) and the game  $v_{\mathbf{x}, \mathcal{D}^{ref}}$  (relative to the distribution of accepted applications  $\mathcal{D}^{accept}$ ). 95% CIs ranged from  $\pm 0.0004$  to  $\pm 0.0007$  for both games.

(estimated using the method described in Section 5.2) show that they are uncertain indeed: the CIs span both positive and negative values.

**Irrelevant references** In Section 4.4, we noted the importance of relevant references (or norms), and how the IME game  $v_{\mathbf{x}}^{unif}$  based on the uniform distribution  $\mathcal{U}$  can focus on irrelevant references. We illustrate this on the Adult Income example; see the third row of Table 4. We find that almost all attribution from the  $v_{\mathbf{x}}^{unif}$  game falls on the *capitalgain* feature. This is surprising as *capitalgain* is zero for the example being explained, and for over 90% of examples in the Adult Income dataset. The attributions are an artifact of uniformly sampling reference feature values, which causes nearly all references to have non-zero capital gain (as the probability of sampling an exactly zero capital gain is infinitesimal).

## 5.2 Applying the FAE framework

We now consider how the FAE framework enables more faithful explanations for the three models we study.

**Formulating contrastive questions** A key benefit of FAE is that it enables explaining predictions relative to a selected group of references. For instance, in the Lending

Club model, rather than asking “Why did our rejected example receive a score of 0.28?” we ask the contrastive question “Why did our rejected example receive a score of 0.28 relative to the examples that were accepted?” This is a more apt question, as it explicitly discards irrelevant comparisons to other rejected applications. In terms of game formulation, the contrastive approach amounts to considering single-reference games where the reference is drawn from the distribution of accepted applications (denoted by  $\mathcal{D}^{accept}$ ) rather than all applications. The attributions for each of these questions (shown in Table 5) turn out to be quite different. For instance, although number of recently-opened accounts (*acc*) is still the highest-attributed feature, we find that credit score (*fico*), income (*inc*), and debt-to-income ratio (*dti*) receive significantly higher attribution in the contrastive formulation. Without formulating the contrastive question, we would be misled into believing that these features are unimportant for the rejection.

**Quantifying uncertainty** When summarizing the attribution distribution with the mean, confidence intervals can be computed using the standard error of the mean (see Section 4.3). Returning to our Bike Sharing example, with 100 samples, the 95% confidence intervals for *hr* and *workingday* are -36 to 15, and -1 to 12, respectively. The large CIs caution us that 100 samples are perhaps too few. When using the full test set, the 95% CIs drop to 0.0 to 5.1 for *hr*, and 0.6 to 2.0 for *workingday*.

**Summarizing attribution distributions** To obtain a more faithful summarization of the single-reference game attributions, we explore a clustering based approach. We compute attributions for single reference games relative to points sampled from the the input distribution ( $\mathcal{D}^{inp}$ ), and then apply *k*-means clustering to these attributions. The resulting clusters effectively group references that yield similar (contrastive) attributions for the prediction at the explained point. Consequently, the attribution distribution within each cluster has a small spread, and can be summarized via the mean.

We applied this approach to all three models and obtained promising results, wherein, clustering helps mine distinct attribution patterns. Table 4 (bottom) shows the results for the Adult Income model; results for other models are provided in Appendix E. Notice that clustering identifies a large group of irrelevant references (cluster 2) which are similar to the explained point, demonstrating low attributions and predictions. Cluster 3 discovers the same pattern that the  $v_x^{unif}$  formulation did: high *capitalgain* causes extremely high scores. Since over 90% of points in the dataset have zero *capitalgain*, this pattern is “washed out in the average” relative to the entire data distribution  $\mathcal{D}^{inp}$  (as in KernelSHAP); see the first row of Table 4. On the other hand, the IME formulation identifies nothing but this pattern. Our clustering also helps identify other patterns. Clusters 1 and 5 show that when compared to references that obtain a high-but-not-extreme score, *marital-status*, *relationship*, and *education-num* are the primary factors accounting for the lower prediction score for the example at hand.

## 6 CONCLUSION

We perform an in-depth study of various Shapley-value-based model explanation methods. We find cases where existing methods yield counter-intuitive attributions, and we

trace these misleading attributions to the cooperative games formulated by these methods. We propose a generalizing formulation that unifies attribution methods, offers clarity for interpreting each method’s attributions, and admits straightforward confidence intervals for attributions.

We propose a conceptual framework for model explanations, called *formulate, approximate, explain* (FAE), which is built on principles from cognitive psychology. We advise practitioners to *formulate* contrastive explanation questions that specify the references relative to which a prediction should be explained, for example “Why did this rejected loan application receive a score of 0.28 *in contrast to the applications that were accepted?*” By *approximating* the Shapley values of games formulated relative to the chosen references, and *explaining* the distribution of approximated Shapley values, we provide a more relevant answer to the explanation question at hand.

Finally, we conclude that axiomatic guarantees do not inherently guarantee relevant explanations, and that game formulations must be constructed carefully. In summarizing attribution distributions, we caution practitioners to avoid coarse-grained summaries, and to quantify any uncertainty resulting from any approximations used.

## References

1. Aas, K., Jullum, M., Løland, A.: Explaining individual predictions when features are dependent: More accurate approximations to shapley values. arXiv preprint arXiv:1903.10464 (2019)
2. Ancona, M., Ceolini, E., Zireli, C., Gross, M.: Towards better understanding of gradient-based attribution methods for deep neural networks. In: International Conference on Learning Representations (2018)
3. Ancona, M., Zireli, C., Gross, M.: Explaining deep neural networks with a polynomial time algorithm for shapley value approximation. In: Proceedings of the 36th International Conference on Machine Learning (2019)
4. Chen, J., Song, L., Wainwright, M.J., Jordan, M.I.: L-shapley and c-shapley: Efficient model interpretation for structured data. arXiv preprint arXiv:1808.02610 (2018)
5. Cohen, S., Ruppin, E., Dror, G.: Feature selection based on the shapley value. In other words **1**, 98Eqr (2005)
6. Datta, A., Sen, S., Zick, Y.: Algorithmic transparency via quantitative input influence: Theory and experiments with learning systems. In: 2016 IEEE symposium on security and privacy (SP). pp. 598–617. IEEE (2016)
7. Dhurandhar, A., Chen, P., Luss, R., Tu, C., Ting, P., Shanmugam, K., Das, P.: Explanations based on the missing: Towards contrastive explanations with pertinent negatives. CoRR (2018), <http://arxiv.org/abs/1802.07623>
8. Efron, B., Tibshirani, R.: Bootstrap methods for standard errors, confidence intervals, and other measures of statistical accuracy. *Statistical Science* pp. 54–75 (1986)
9. Ghorbani, A., Zou, J.: Data shapley: Equitable valuation of data for machine learning. In: Proceedings of the 36th International Conference on Machine Learning (2019)
10. Hesslow, G.: The problem of causal selection. In: Hilton, D.J. (ed.) *Contemporary Science and Natural Explanation: Commonsense Conceptions of Causality*. New York University Press (1988)
11. Hitchcock, C., Knobe, J.: Cause and norm. *Journal of Philosophy* **106**(11), 587–612 (2009)

12. Holzinger, A., Kickmeier-Rust, M., Mller, H.: KANDINSKY Patterns as IQ-Test for Machine Learning. In: Holzinger, A., Kieseberg, P., Tjoa, A.M., Weippl, E. (eds.) *Machine Learning and Knowledge Extraction*. pp. 1–14. Lecture Notes in Computer Science, Springer International Publishing, Cham (2019).
13. Hunt, X.J., Abbey, R., Tharrington, R., Huiskens, J., Wesdorp, N.: An ai-augmented lesion detection framework for liver metastases with model interpretability. *arXiv preprint arXiv:1907.07713* (2019)
14. Janzing, D., Minorics, L., Blbaum, P.: Feature relevance quantification in explainable ai: A causal problem. *arXiv preprint arXiv:1910.13413* (2019)
15. Kahneman, D., Miller, D.T.: Norm theory: Comparing reality to its alternatives. *Psychological review* **93**(2), 136 (1986)
16. Ke, G., Meng, Q., Finley, T., Wang, T., Chen, W., Ma, W., Ye, Q., Liu, T.Y.: Lightgbm: A highly efficient gradient boosting decision tree. In: *Advances in Neural Information Processing Systems*. pp. 3146–3154 (2017)
17. Lipton, P.: Contrastive explanation. *Royal Institute of Philosophy Supplement* **27**, 247266 (1990). <https://doi.org/10.1017/S1358246100005130>
18. Lundberg, S.M., Erion, G.G., Lee, S.I.: Consistent individualized feature attribution for tree ensembles. *arXiv preprint arXiv:1802.03888* (2018)
19. Lundberg, S.M., Lee, S.I.: A unified approach to interpreting model predictions. In: *Advances in Neural Information Processing Systems*. pp. 4765–4774 (2017)
20. Maleki, S., Tran-Thanh, L., Hines, G., Rahwan, T., Rogers, A.: Bounding the estimation error of sampling-based shapley value approximation. *arXiv preprint arXiv:1306.4265* (2013)
21. Miller, T.: Explanation in artificial intelligence: Insights from the social sciences. *arXiv preprint arXiv:1706.07269* (2017)
22. Mittelstadt, B., Russell, C., Wachter, S.: Explaining explanations in ai. In: *Proceedings of the conference on fairness, accountability, and transparency*. pp. 279–288. ACM (2019)
23. Ribeiro, M.T., Singh, S., Guestrin, C.: Why should i trust you?: Explaining the predictions of any classifier. In: *SIGKDD international conference on knowledge discovery and data mining*. pp. 1135–1144 (2016)
24. Shapley, L.S.: A value for n-person games. *Contributions to the Theory of Games* **2**(28), 307–317 (1953)
25. Shrikumar, A., Greenside, P., Kundaje, A.: Learning important features through propagating activation differences. In: *34th International Conference on Machine Learning-Volume 70*. pp. 3145–3153 (2017)
26. Štrumbelj, E., Kononenko, I.: An efficient explanation of individual classifications using game theory. *Journal of Machine Learning Research* **11**, 1–18 (2010)
27. Štrumbelj, E., Kononenko, I.: Explaining prediction models and individual predictions with feature contributions. *Knowledge and information systems* **41**(3), 647–665 (2014)
28. Sundararajan, M., Najmi, A.: The many shapley values for model explanation. *arXiv preprint arXiv:1908.08474* (2019)
29. Sundararajan, M., Taly, A., Yan, Q.: Axiomatic attribution for deep networks. In: *Proceedings of the 34th International Conference on Machine Learning-Volume 70*. pp. 3319–3328. JMLR. org (2017)
30. Young, H.P.: Monotonic solutions of cooperative games. *International Journal of Game Theory* **14**, 65–72 (1985)

## Appendix

### A Shapley Value Axioms

We briefly summarize the four Shapley value axioms.

- The *Dummy* axiom requires that if player  $i$  has no possible contribution (i.e.  $v(S \cup \{i\}) = v(S)$  for all  $S \subseteq \mathcal{M}$ ), then that player receives zero attribution.
- The *Symmetry* axiom requires that two players that always have the same contribution receive equal attribution, Formally, if  $v(S \cup \{i\}) = v(S \cup \{j\})$  for all  $S$  not containing  $i$  or  $j$  then  $\phi_i(v) = \phi_j(v)$ .
- The *Efficiency* axiom requires that the attributions to all players sum to the total payoff of all players. Formally,  $\sum_i \phi_i(v) = v(\mathcal{M})$ .
- The *Linearity* axiom states that for any payoff function  $v$  that is a linear combination of two other payoff functions  $u$  and  $w$  (i.e.  $v(S) = \alpha u(S) + \beta w(S)$ ), the Shapley values of  $v$  equal the corresponding linear combination of the Shapley values of  $u$  and  $w$  (i.e.  $\phi_i(v) = \alpha \phi_i(u) + \beta \phi_i(w)$ ).

### B Additional Shapley value approximations

**Marginal contribution sampling** We can express the Shapley value of a player as the expected value of the weighted marginal contribution to a random coalition  $S$  sampled uniformly from all possible coalitions excluding that player, rather than an exhaustive weighted sum. A sampling estimator of this expectation is by nature unbiased, so this can be used as an alternative to the permutation estimator in approximating attributions with confidence intervals.

$$\phi_i(v) = \mathbb{E}_S \left[ \frac{2^{M-1}}{M} \binom{M-1}{|S|}^{-1} (v(S \cup \{i\}) - v(S)) \right] \quad (13)$$

Equation 13 can be approximated with a Monte Carlo estimate, i.e. by sampling from the random  $S$  and averaging the quantity within the expectation.

**Weighted least squares** The Shapley values are the solution to a certain weighted least squares optimization problem which was popularized through its use in the KernelSHAP algorithm. For a full explanation, see <https://arxiv.org/abs/1903.10464>.

$$\phi = \arg \min_{\phi} \sum_{S \subseteq \mathcal{M}} \frac{M-1}{\binom{M}{|S|} |S| (M-|S|)} \left( v(S) - \sum_{i=1}^M \phi_i \right)^2 \quad (14)$$

The fraction in the left of Equation 14 is often referred to as the Shapley kernel. In practice, an approximate objective function is minimized. The approximate objective is defined as a summation over squared error on a sample of coalitions rather than over squared error on all possible coalitions. Additionally, the “KernelSHAP trick” may be employed, wherein sampling is performed according to the Shapley kernel (rather than

uniformly), and the least-squares optimization is solved with uniform weights (rather than Shapley kernel weights) to account for the adjusted sampling.

To the best of our knowledge, there exists no proof that the solution to a subsampled objective function of the form in Equation 14 is an estimator (unbiased or otherwise) of the Shapley values. In practice, it does appear that subsampling down to even a small fraction of the total number of possible coalitions (weighted by the Shapley kernel or uniformly) does a good job of estimating the Shapley values for explanation games. Furthermore, approximation errors in such experiments do not yield signs of bias. However, we do note that using the weighted least squares approximation with our confidence interval equation does inherently imply an unproved assumption that it is an unbiased estimator.

## C Proofs

In what follows, we prove the lemmas from the main paper. The proofs refer to equations and definition from the main paper.

### C.1 Proof of Proposition 1

From the definitions of  $v_{\mathbf{x}, \mathcal{D}^{ref}}$  (Equation 8) and  $v_{\mathbf{x}, \mathbf{r}}$  (Equation 9), it follows that  $v_{\mathbf{x}, \mathcal{D}^{ref}}(S) = \mathbb{E}_{\mathbf{R} \sim \mathcal{D}^{ref}} [v_{\mathbf{x}, \mathbf{R}}(S)]$ . Thus, the game  $v_{\mathbf{x}, \mathcal{D}^{ref}}$  is a linear combination of games  $\{v_{\mathbf{x}, \mathbf{r}} \mid \mathbf{r} \in \mathcal{X}\}$  with weights defined by the distribution  $\mathcal{D}^{ref}$ . From the Linearity axiom of Shapley values, it follows that the Shapley values of the game  $v_{\mathbf{x}, \mathcal{D}^{ref}}$  must be a corresponding linear combination of the Shapley values of the games  $\{v_{\mathbf{x}, \mathbf{r}} \mid \mathbf{r} \in \mathcal{X}\}$  (with weights defined by the distribution  $\mathcal{D}^{ref}$ ). Thus,  $\phi(v_{\mathbf{x}, \mathcal{D}^{ref}}) = \mathbb{E}_{\mathbf{R} \sim \mathcal{D}^{ref}} [\phi(v_{\mathbf{x}, \mathbf{R}})]$ .  $\square$

### C.2 Proof of Proposition 2

From Proposition 1, we have  $\phi_i(v_{\mathbf{x}, \mathcal{D}^{ref}}) = \mathbb{E}_{\mathbf{R} \sim \mathcal{D}^{ref}} [\phi_i(v_{\mathbf{x}, \mathbf{R}})]$ . Thus, to prove this lemma, it suffices to show that for any irrelevant feature  $i$ , the Shapley value from the game  $v_{\mathbf{x}, \mathbf{r}}$  is zero for all references  $\mathbf{r} \in \mathcal{X}$ . That is,

$$\forall \mathbf{r} \in \mathcal{X} \phi_i(v_{\mathbf{x}, \mathbf{r}}) = 0 \quad (15)$$

From the definition of Shapley values (Equation 1), we have:

$$\phi_i(v_{\mathbf{x}, \mathbf{r}}) = \frac{1}{M} \sum_{S \subseteq \mathcal{M} \setminus \{i\}} \binom{M-1}{|S|}^{-1} (v_{\mathbf{x}, \mathbf{r}}(S \cup \{i\}) - v_{\mathbf{x}, \mathbf{r}}(S)) \quad (16)$$

Thus, to prove Equation 15 it suffices to show the marginal contribution  $(v_{\mathbf{x}, \mathbf{r}}(S \cup \{i\}) - v_{\mathbf{x}, \mathbf{r}}(S))$  of an irrelevant feature  $i$  to any subset of features  $S \subseteq \mathcal{M} \setminus \{i\}$  is always zero. From the definition of the game  $v_{\mathbf{x}, \mathbf{r}}$ , we have:

$$v_{\mathbf{x}, \mathbf{r}}(S \cup \{i\}) - v_{\mathbf{x}, \mathbf{r}}(S) = f(\mathbf{z}(\mathbf{x}, \mathbf{r}, S \cup \{i\})) - f(\mathbf{z}(\mathbf{x}, \mathbf{r}, S)) \quad (17)$$

From the definition of composite inputs  $\mathbf{z}$  (Equation 3), it follows that the inputs  $\mathbf{z}(\mathbf{x}, \mathbf{r}, S \cup \{i\})$  and  $\mathbf{z}(\mathbf{x}, \mathbf{r}, S)$  agree on all features except  $i$ . Thus, if feature  $i$  is irrelevant,  $f(\mathbf{z}(\mathbf{x}, \mathbf{r}, S \cup \{i\})) = f(\mathbf{z}(\mathbf{x}, \mathbf{r}, S))$ , and consequently by Equation 16,  $v_{\mathbf{x}, \mathbf{r}}(S \cup \{i\}) - v_{\mathbf{x}, \mathbf{r}}(S) = 0$  for all subsets  $S \subseteq \mathcal{M} \setminus \{i\}$ . Combining this with the definition of Shapley values (Equation 1) proves Equation 15.  $\square$

## D Reproducibility

For brevity, we omitted from the main paper many of the mundane choices in the design of our toy examples and case studies. To further transparency and reproducibility, we include them here.

### D.1 Fitting models

For both case studies, we used the LightGBM package configured with default parameters to fit a Gradient Boosted Decision Trees (GBDT) model. For the Bike Sharing dataset, we fit on all examples from 2011 while holding out the 2012 examples for testing. We omitted the *atemp* feature, as it is highly correlated to *temp* ( $r = 0.98$ ), and the *instant* feature because the tree-based GBDT model cannot capture its time-series trend. For parsimony, we refitted the model to the top five most important features by cumulative gain (*hr*, *temp*, *workingday*, *hum*, and *season*). This lowered test-set  $r^2$  from 0.64 to 0.63. For the Adult Income dataset, we used the pre-defined train/test split. Again, we refitted the model to the top five features by cumulative gain feature importance (*relationship*, *capitalgain*, *education-num*, *marital-status*, and *age*). This increased test-set misclassification error from 14.73% to 10.97%.

### D.2 Selection of points to explain

For the Bike Share case study, we sampled ten points at random from the test set. We selected one whose prediction was close to the middle of the range observed over the entire test set (predictions ranged approximately from 0 to 600). Specifically, we selected instant 11729 (2012-05-08, 9pm). We examined other points from the same sample of ten to suggest a random but meaningful comparative question. We found another point with comparable *workingday*, *hum*, and *season*: instant 11362. This point caught our eye because it differed only in *hr* (2pm rather than 9pm), and *temp* (0.36 rather than 0.64) but had a much lower prediction.

For the Adult Income case study, we wanted to explain why a point was scored as likely to have low income, a task roughly analogous to that of explaining why an application for credit is rejected by a creditworthiness model in a lending setting. We sampled points at random with scores between 0.01 and 0.1, and chose the 9880th point in the test set due to its strikingly high *education-num* (most of the low-scoring points sampled had lower *education-num*).

For the Lending Club data, we chose an open-source subset of the dataset that has been pre-cleaned to a predictive task on 3-year loans. For the five-feature model, we selected the top five features by cumulative gain feature importance from a model fit to the full set of features.

### D.3 K-means clustering

We choose  $k = 5$  arbitrarily, having observed a general tradeoff of conciseness for precision as  $k$  increases. In the extremes,  $k = 1$  maintains the overall attribution distribution, while  $k = N$  examines each single-reference game separately.

## E Case Study Supplemental Material

Game Formulation	Size	Avg. Prediction ( $\phi_0$ )	hr	temp	work.	hum	season
$v_{\mathbf{x}}^{inp}$	100%	151	3	47	1	7	2
$v_{\mathbf{x}}^{J.M.}$	100%	141	6	50	1	9	3
$v_{\mathbf{x}}^{unif}$	100%	128	3	60	3	12	3
Cluster 1	12.9%	309	-86	14	-28	3	-1
Cluster 2	27.6%	28	140	32	0	9	0
Cluster 3	10.5%	375	-247	58	16	9	-1
Cluster 4	32.5%	131	31	38	3	4	2
Cluster 5	16.5%	128	-57	107	13	9	9

Table 6: Bike Sharing comparison of mean attributions. 95% CIs ranged from  $\pm 0.4$  (*hum* in  $\mathcal{D}^{inp}$  and  $\mathcal{D}^{J.M.}$ ) to  $\pm 2.5$  (*hr* in  $\mathcal{D}^{inp}$  and  $\mathcal{D}^{J.M.}$ ).

Game	Size	Avg. Prediction ( $\phi_0$ )	rel.	cap.	edu.	mar.	age
$v_{\mathbf{x}}^{inp}$	100%	0.24	-0.04	-0.03	-0.01	-0.10	-0.00
$v_{\mathbf{x}}^{J.M.}$	100%	0.19	-0.02	-0.03	-0.01	-0.08	0.01
$v_{\mathbf{x}}^{unif}$	100%	0.82	0.01	-0.79	0.02	-0.03	0.04
Cluster 1	10.2%	0.67	-0.15	-0.01	-0.15	-0.28	-0.02
Cluster 2	55.3%	0.04	0.01	0.00	0.00	-0.01	0.02
Cluster 3	4.4%	0.99	-0.04	-0.70	-0.06	-0.12	-0.01
Cluster 4	28.0%	0.31	-0.09	0.00	0.08	-0.21	-0.03
Cluster 5	2.1%	0.67	-0.04	0.01	-0.47	-0.14	0.03

Table 7: Adult Income comparison of mean attributions. 95% CIs ranged from  $\pm 0.0004$  (Cluster 2, *relationship*) to  $\pm 0.0115$  (Cluster 5, *marital-status* and *age*).

Game	Size	Avg. Prediction ( $\phi_0$ )	fico.	addr.	inc.	acc.	dti
$v_{\mathbf{x}, \mathcal{D}^{ref}}$	20%	0.05	0.02	0.04	0.02	0.11	0.03
$v_{\mathbf{x}}^{inp}$	100%	0.14	0.00	0.03	0.00	0.10	0.00
$v_{\mathbf{x}}^{J.M.}$	100%	0.14	0.01	0.03	0.01	0.10	0.00
$v_{\mathbf{x}}^{unif}$	100%	0.11	0.05	0.07	-0.01	0.03	0.02
Cluster 1	28.5%	0.11	0.01	0.06	0.00	0.08	0.01
Cluster 2	24.4%	0.10	0.01	0.00	0.01	0.11	0.04
Cluster 3	15.4%	0.18	0.00	0.01	0.00	0.14	-0.05
Cluster 4	17.6%	0.16	-0.01	0.01	0.03	0.09	-0.01
Cluster 5	14.0%	0.22	-0.01	0.05	-0.02	0.08	-0.06

Table 8: Lending Club comparison of mean attributions. 95% CIs ranged from  $\pm 0.0004$  to  $\pm 0.0007$  for all games.