

Approximating the Shapley Value without Marginal Contributions

Patrick Kolpaczki¹, Viktor Bengs^{2,3}, Maximilian Muschalik^{2,3}, Eyke Hüllermeier^{2,3}

¹Paderborn University

²Institute of Informatics, University of Munich (LMU)

³Munich Center for Machine Learning

patrick.kolpaczki@upb.de, viktor.bengs@lmu.de, maximilian.muschalik@lmu.de, eyke@lmu.de

Abstract

The Shapley value, which is arguably the most popular approach for assigning a meaningful contribution value to players in a cooperative game, has recently been used intensively in explainable artificial intelligence. Its meaningfulness is due to axiomatic properties that only the Shapley value satisfies, which, however, comes at the expense of an exact computation growing exponentially with the number of agents. Accordingly, a number of works are devoted to the efficient approximation of the Shapley value, most of them revolve around the notion of an agent’s marginal contribution. In this paper, we propose with *SVARM* and *Stratified SVARM* two parameter-free and domain-independent approximation algorithms based on a representation of the Shapley value detached from the notion of marginal contribution. We prove unmatched theoretical guarantees regarding their approximation quality and provide empirical results including synthetic games as well as common explainability use cases comparing ourselves with state-of-the-art methods.

Introduction

Whenever agents can federalize in groups (form coalitions) to accomplish a task and get rewarded with a collective benefit that is to be shared among the group members, the notion of *cooperative game* stemming from game theory is arguably the most favorable concept to model such situations. This is due to its simplicity, which nevertheless allows for covering a whole range of practical applications. The agents are called *players* and are contained in a player set \mathcal{N} . Each possible subset of players $S \subseteq \mathcal{N}$ is understood as a *coalition* and the coalition \mathcal{N} containing all players is called the *grand coalition*. The collective benefit $\nu(S)$ that a coalition S receives upon formation is given by a *value function* ν assigning each coalition a real-valued *worth*.

The connection of cooperative games to (supervised) machine learning is already well-established. The most prominent example is feature importance scores, both local and global, for a machine learning model: features of a dataset can be seen as players, allowing one to interpret a feature subset as a coalition, while the model’s generalization performance using exactly that feature subset is its worth (Cohen, Dror, and Ruppin 2007). Other applications include

evaluating the importance of parameters in a machine learning model, e.g. single neurons in a deep neural network (Ghorbani and Zou 2020) or base learners in an ensemble (Rozenberczki and Sarkar 2021), or assigning relevance scores to datapoints in a given dataset (Ghorbani and Zou 2019). See Rozenberczki et al. (2022) for a wider overview of its usage in the field of explainable artificial intelligence. Outside the realm of machine learning cooperative games also found applications in operations research (Luo, Zhou, and Lev 2022), for finding fair compensation mechanisms in electricity grids (O’Brien, Gamal, and Rajagopal 2015), or even for the purpose of identifying the most influential individuals in terrorist networks (van Campen et al. 2018).

In all of these applications, the question naturally arises of how to appropriately determine the contribution of a single player (feature, parameter, etc.) with respect to the grand collective benefit. In other words, how to allocate the worth $\nu(\mathcal{N})$ of the full player set \mathcal{N} among the players in a fair manner. The indisputably most popular solution to this problem is the *Shapley value* (Shapley 1953), which can be intuitively expressed by *marginal contributions*. We call the increase in worth that comes with the inclusion of player i to a coalition S , i.e., the difference $\nu(S \cup \{i\}) - \nu(S)$, player i ’s marginal contribution to S . The Shapley value of i is a weighted average of all its marginal contributions to coalitions that do not include i . Its popularity stems from the fact that it is the only solution to satisfy axiomatic properties that arguably capture fairness (Shapley 1953).

Despite the appealing theoretical properties of the Shapley value, there is one major drawback with respect to its practical application, as its computational complexity increases exponentially with the number of players n . As a consequence, the exact computation of the Shapley value becomes practically infeasible even for a moderate number of players. This is especially the case where accesses to ν are costly, e.g., re-evaluating a (complex) machine learning model for a specific feature subset, or manipulating training data each time ν is accessed. Recently, several approximation methods have been proposed in search of a remedy, enabling the utilization of the Shapley value in explainable AI (and beyond). However, most works are stiffened towards the notion of marginal contribution, and, consequently, judge algorithms by their achieved approximation accuracy depending on the number of evaluated marginal

contributions. This measure does not do justice to the fact that approximations can completely dispense with the consideration of marginal contributions and elicit information from ν in a more efficient way — as we show in this paper. We claim that the number of single accesses to ν should be considered instead, since especially in machine learning, as mentioned above, access to ν is a bottleneck in overall runtime. In this paper, we make up for this deficit by considering the problem of approximating the Shapley values under a fixed *budget* T of evaluations (accesses) of ν .

Contribution. We present a novel representation of the Shapley value that does not rely on the notion of marginal contribution. Our first proposed approximation algorithm *Shapley Value Approximation without Requesting Marginals* (SVARM) exploits this representation and directly samples values of coalitions, facilitating “a swarm of updates”, i.e., multiple Shapley value estimates are updated at once. This is in stark contrast to the usual way of sampling marginal contributions that only allows the update of a single estimate. We prove theoretical guarantees regarding SVARM’s precision including the bound of $\mathcal{O}(\frac{\log n}{T-n})$ on its variance.

Based on a partitioning of the set of all coalitions according to their size, we develop with *Stratified SVARM* a refinement of SVARM. The applied stratification materializes a twofold improvement: (i) the homogeneous strata (w.r.t. the coalition worth) significantly accelerate convergence of estimates, (ii) our stratified representation of the Shapley value with decomposed marginal contributions facilitates a mechanism that updates the estimates of *all* players with *each single* coalition sampled. Among other results, we bound its variance by $\mathcal{O}(\frac{\log n}{T-n \log n})$.

Besides our superior theoretical findings, both algorithms possess a number of properties in their favor. More specifically, both are unbiased, parameter-free, incremental, i.e., the available budget has not to be fixed and can be enlarged or cut prematurely, facilitating on-the-fly approximations due to their anytime property, and do not require any knowledge about the latent value function. Moreover, both are domain-independent and not limited to some specific fields, but can be used to approximate the Shapley values of any possible cooperative game.

Finally, we compare our algorithms empirically against other popular competitors, demonstrating their practical usefulness and proving our empirical enhancement *Stratified SVARM*⁺, which samples without replacement to be the first sample-mean-based approach to achieve rivaling state-of-the-art approximation quality. All code including documentation and the technical appendix can be found on GitHub¹.

Related Work

The recent rise of explainable AI has incentivized the research on approximation methods for the Shapley value leading to a variety of different algorithms for this purpose. The first distinction to be made is between those that are domain-independent, i.e., able to deal with any cooperative

game, and those that are tailored to a specific use case, e.g. assigning Shapley values to single neurons in neural networks, or which impose specific assumptions on the value function. In this paper, we will consider only the former, as it is our goal to provide approximation algorithms independent of the context in which they are applied. The first and so far simplest of this kind is *ApproShapley* (Castro, Gómez, and Tejada 2009), which samples marginal contributions from each player based on randomly drawn permutations of the player set. The variance of each of its Shapley value estimates is bounded by $\mathcal{O}(\frac{n}{T})$. *Stratified Sampling* (Maleki et al. 2013) and *Structured Sampling* (van Campen et al. 2018) both partition the marginal contributions of each player by coalition size in order to stratify the marginal contributions of the population from which to draw a sample, which leads to a variance reduction. While *Stratified Sampling* calculates a sophisticated allocation of samples for each coalition size, *Structured Sampling* simply samples with equal frequencies. Multiple follow-up works suggest specific techniques to improve the sampling allocation over the different coalition sizes (O’Brien, Gamal, and Rajagopal 2015; Castro et al. 2017; Burgess and Chapman 2021).

In order to reduce the variance of the naive sampling approach underlying *ApproShapley*, Illés and Kerényi (2019) suggest to use ergodic sampling, i.e., generating samples that are not independent but still satisfy the strong Law of Large numbers. Quite recently, Mitchell et al. (2022) investigated two techniques for improving *ApproShapley*’s sampling approach. One is based on the theory of reproducing kernel Hilbert spaces, which focuses on minimizing the discrepancies for functions of permutations. The other exploits a geometrical connection between uniform sampling on the Euclidean sphere and uniform sampling over permutations.

Adopting a Bayesian perspective, i.e., by viewing the Shapley values as random variables, Touati, Radjef, and Sais (2021) consider approximating the Shapley values by Bayesian estimates (posterior mean, mode, or median), where each posterior distribution of a player’s Shapley value depends on the remaining ones. Utilizing a representation of the Shapley value as an integral (Owen 1972), *Owen Sampling* (Okhrati and Lipani 2020) approximates this integral by sampling marginal contributions using antithetic sampling (Rubinstein and Kroese 2016; Lomeli et al. 2019) for variance reduction.

A fairly new class of approaches that dissociates itself from the notion of marginal contribution are those that view the Shapley value as a solution of a quadratic program with equality constraints (Lundberg and Lee 2017; Simon and Vincent 2020; Covert and Lee 2021). Another unorthodox approach is to divide the player set into small enough groups for which the Shapley values within these groups can be computed exactly (Soufiani et al. 2014; Corder and Decker 2019). For an overview of approaches related to machine learning we refer to (Chen et al. 2023).

Problem Statement

The formal notion of a cooperative game is defined by a tuple (\mathcal{N}, ν) consisting of a set of players $\mathcal{N} = \{1, \dots, n\}$ and a value function $\nu : \mathcal{P}(\mathcal{N}) \rightarrow \mathbb{R}$ that assigns to each subset

¹<https://github.com/kolpaczki/Approximating-the-Shapley-Value-without-Marginal-Contributions>

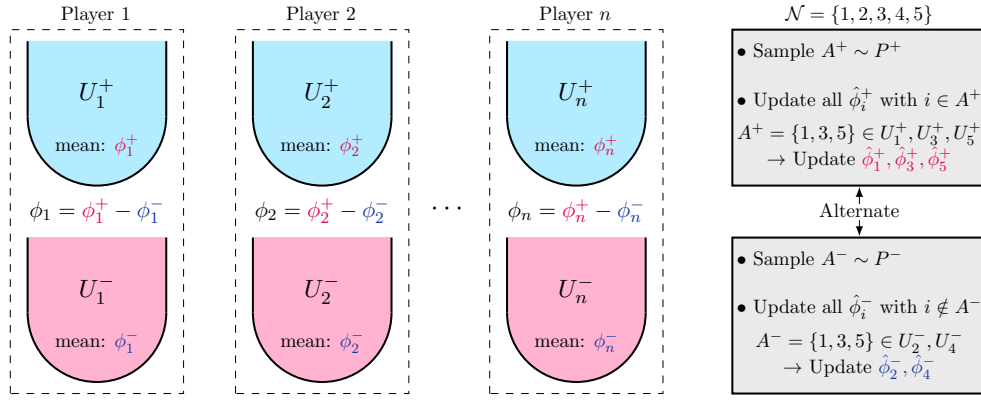


Figure 1: Illustration of SVARM’s sampling process and update rule: Each player i has two urns $U_i^+ := \{S \cup \{i\} \mid S \subseteq \mathcal{N}_i\}$ and $U_i^- := \{S \mid S \subseteq \mathcal{N}_i\}$ containing marbles which represent coalitions, with mean coalition worth ϕ_i^+ and ϕ_i^- . SVARM alternates between sampling coalitions $A^+ \sim P^+$ and $A^- \sim P^-$. With each drawn coalition all estimates of those urns are updated which contain the corresponding marble. Since each player’s two urns form a partition of the powerset $\mathcal{P}(\mathcal{N})$, all players have exactly one urn updated with each sample.

of \mathcal{N} a real-valued number. The value function must satisfy $\nu(\emptyset) = 0$. We call the subsets of \mathcal{N} coalitions, \mathcal{N} itself the grand coalition, and the assigned value $\nu(S)$ to a coalition $S \subseteq \mathcal{N}$ its worth. Given a cooperative game (\mathcal{N}, ν) , the Shapley value assigns each player a share of the grand coalition’s worth. In particular, the Shapley value (Shapley 1953) of any player $i \in \mathcal{N}$ is defined as

$$\phi_i = \sum_{S \subseteq \mathcal{N}_i} \frac{1}{n \cdot \binom{n-1}{|S|}} [\nu(S \cup \{i\}) - \nu(S)], \quad (1)$$

where $\mathcal{N}_i := \mathcal{N} \setminus \{i\}$ for each player $i \in \mathcal{N}$. The term $\nu(S \cup \{i\}) - \nu(S)$ is also known as player i ’s marginal contribution to $S \subseteq \mathcal{N}_i$ and captures the increase in collective benefit when player i joins the coalition S . Thus, the Shapley value can be seen as the weighted average of a player’s marginal contributions.

The exact computation of all Shapley values requires the knowledge of the values of all 2^n many coalitions² and is shown to be NP-hard (Deng and Papadimitriou 1994). In light of the exponential computational effort w.r.t. to n , we consider the goal of approximating the Shapley value of all players as precisely as possible for a given *budget* of $T \in \mathbb{N}$ many evaluations (accesses) of ν in discrete time steps $1, \dots, T$. Since $\nu(\emptyset) = 0$ holds by definition, the evaluation of $\nu(\emptyset)$ comes for free without any budget cost. We judge the quality of the estimates $\hat{\phi}_1, \dots, \hat{\phi}_n$ — which are possibly of stochastic nature — obtained by an approximation algorithm after T many evaluations by two criteria that have to be minimized for all $i \in \mathcal{N}$. First, the mean squared error (MSE) of the estimate $\hat{\phi}_i$ is given by

$$\mathbb{E}[(\hat{\phi}_i - \phi_i)^2]. \quad (2)$$

Utilizing the bias-variance decomposition allows us to reduce the squared error to the variance $\mathbb{V}[\hat{\phi}_i]$ of the Shapley

²In fact, only $2^n - 1$ many coalitions, as $\nu(\emptyset) = 0$ is known.

value estimate in case that it is unbiased, i.e. $\mathbb{E}[\hat{\phi}_i] = \phi_i$. The second criterion is the probability of $\hat{\phi}_i$ deviating from ϕ_i by more than a fixed $\varepsilon > 0$:

$$\mathbb{P}(|\hat{\phi}_i - \phi_i| > \varepsilon). \quad (3)$$

Both criteria are well-established for measuring the quality of an algorithm approximating the Shapley value.

SVARM

Thanks to the distributive law, the formula of the Shapley value for a player i can be rearranged so that it is not its weighted average of marginal contributions, but the difference of the weighted average of coalition values by adding i and the weighted average of coalition values without i :

$$\phi_i = \underbrace{\sum_{S \subseteq \mathcal{N}_i} w_S \cdot \nu(S \cup \{i\})}_{=: \phi_i^+} - \underbrace{\sum_{S \subseteq \mathcal{N}_i} w_S \cdot \nu(S)}_{=: \phi_i^-}, \quad (4)$$

with weights $w_S = \frac{1}{n \cdot \binom{n-1}{|S|}}$ for each $S \subseteq \mathcal{N}_i$. We call ϕ_i^+ the positive and ϕ_i^- the negative Shapley value, while we refer to the collective of both as the signed Shapley values. The weighted averages ϕ_i^+ and ϕ_i^- can also be viewed as expected values, i.e., $\phi_i^+ = \mathbb{E}[\nu(S \cup \{i\})]$ and $\phi_i^- = \mathbb{E}[\nu(S)]$, where $S \sim P^w$ and $P^w(S) = w_S$ for all $S \subseteq \mathcal{N}_i$. Note that all weights add up to 1 and thus P^w forms a well-defined probability distribution. In this way, we can approximate each signed Shapley value separately using estimates $\hat{\phi}_i^+$ and $\hat{\phi}_i^-$ and combine them into a Shapley value estimate by means of $\hat{\phi}_i = \hat{\phi}_i^+ - \hat{\phi}_i^-$.

In light of this, a naive approach for approximating each signed Shapley value of a player is by sampling some number of M many coalitions $S^{(1)}, \dots, S^{(M)}$ with distribution P^w and using the sample mean as the estimate, i.e., $\hat{\phi}_i^+ = \frac{1}{M} \sum_{m=1}^M \nu(S^{(m)} \cup \{i\})$. However, this would require all $2n$ signed Shapley values (two per player) to be

estimated separately by sampling coalitions in a dedicated manner, each of which would lead to an update of only one estimate. This ultimately slows down the convergence of the estimates, especially for large n .

On the basis of the aforementioned representation of the Shapley value, we present the *Shapley Value Approximation without Requesting Marginals* (SVARM) algorithm, a novel approach that updates multiple Shapley value estimates at once with a single evaluation of ν . Its novelty consists of sampling coalitions independently from two specifically chosen distributions P^+ and P^- in an alternating fashion, which allows for a more powerful update rule: each (independently) sampled coalition A^+ from P^+ allows one to update all positive Shapley value estimates $\hat{\phi}_i^+$ of all payers i which are contained in A^+ , i.e., $i \in A^+$. Likewise, for a coalition A^- drawn from P^- , all negative Shapley value estimates $\hat{\phi}_i^-$ for $i \notin A^-$ can be updated.

It is worth noting that, for simplicity, we alternate evenly between the samples from the P^+ and P^- distributions, although one could also use a ratio other than $1/2$. To avoid a bias, both distributions have to be tailored such that the following holds for all $i \in \mathcal{N}$ and $S \subseteq \mathcal{N}_i$:

$$\begin{aligned} \mathbb{P}(A^+ = S \cup \{i\} \mid i \in A^+) &= \\ &= \mathbb{P}(A^- = S \mid i \notin A^-) = w_S. \end{aligned} \quad (5)$$

For this reason, we define the probability distributions over coalitions to sample from as

$$P^+(S) := \frac{1}{|S| \binom{n}{|S|} H_n} \quad \forall S \in \mathcal{P}(\mathcal{N}) \setminus \{\emptyset\}, \quad (6)$$

$$P^-(S) := \frac{1}{(n - |S|) \binom{n}{|S|} H_n} \quad \forall S \in \mathcal{P}(\mathcal{N}) \setminus \{\mathcal{N}\}, \quad (7)$$

where $H_n = \sum_{k=1}^n 1/k$ denotes the n -th harmonic number. Note that both P^+ and P^- assign equal probabilities to coalitions of the same size, so that one can first sample the size and then draw a set uniformly of that size. This pair of distributions is probably the only one to fulfill the required property (see Appendix C.1).

The approach of dividing the Shapley value into two parts and approximating both has already been pursued (although not as formally rigorous) via importance sampling (Covert, Lundberg, and Lee 2019), allowing to update all n estimates with each sample. Wang and Jia (2023) adopt the same representation for the Banzhaf value, and coined the strategy of updating all players' estimates with each sampled coalition the *maximum sample reuse* (MSR) principle. Their approximation algorithm is specifically tailored to the Banzhaf value as it leverages its uniform weights $w_S = \frac{1}{2^{n-1}}$ and is thus, at least not directly, transferable to the Shapley value.

In the following we describe SVARM's procedure with the pseudocode of Algorithm 1. The overall idea of the sampling and update process is illustrated in Figure 1. It starts by initializing the positive and negative Shapley value estimates $\hat{\phi}_i^+$ and $\hat{\phi}_i^-$, and the number of samples c_i^+ and c_i^- collected for each player i . SVARM continues by launching a warm-up phase (see Algorithm 3 in Appendix B). In the main loop, the update rule is applied for as many sampled

Algorithm 1: SVARM

Input: $\mathcal{N}, T \in \mathbb{N}$

```

1:  $\hat{\phi}_i^+, \hat{\phi}_i^- \leftarrow 0$  for all  $i \in \mathcal{N}$ 
2:  $c_i^+, c_i^- \leftarrow 1$  for all  $i \in \mathcal{N}$ 
3: WARMUP
4:  $t \leftarrow 2n$ 
5: while  $t + 2 \leq T$  do
6:   Draw  $A^+ \sim P^+$ 
7:   Draw  $A^- \sim P^-$ 
8:    $v^+ \leftarrow \nu(A^+)$ 
9:    $v^- \leftarrow \nu(A^-)$ 
10:  for  $i \in A^+$  do
11:     $\hat{\phi}_i^+ \leftarrow \frac{c_i^+ \hat{\phi}_i^+ + v^+}{c_i^+ + 1}$ 
12:     $c_i^+ \leftarrow c_i^+ + 1$ 
13:  end for
14:  for  $i \in \mathcal{N} \setminus A^-$  do
15:     $\hat{\phi}_i^- \leftarrow \frac{c_i^- \hat{\phi}_i^- + v^-}{c_i^- + 1}$ 
16:     $c_i^- \leftarrow c_i^- + 1$ 
17:  end for
18:   $t \leftarrow t + 2$ 
19: end while
20:  $\hat{\phi}_i \leftarrow \hat{\phi}_i^+ - \hat{\phi}_i^-$  for all  $i \in \mathcal{N}$ 

```

Output: $\hat{\phi}_1, \dots, \hat{\phi}_n$

pairs of coalitions A^+ and A^- as possible until SVARM runs out of budget. In each iteration A^+ is sampled from P^+ and A^- from P^- . The worth of A^+ and A^- is evaluated and stored in v^+ and v^- , requiring two accesses to the value function. The estimate $\hat{\phi}_i^+$ of each player $i \in A^+$ is updated with the worth $\nu(A^+)$ such that $\hat{\phi}_i^+$ is the mean of sampled coalition values. Likewise, the estimate $\hat{\phi}_i^-$ of each player $i \notin A^-$ is updated with the worth $\nu(A^-)$. At the same time, the sample numbers of the respective signed Shapley value estimates are also updated. Finally, SVARM computes its Shapley value estimate $\hat{\phi}_i$ of ϕ_i for each i according to Equation (4). Note that since only the quantities $\hat{\phi}_i^+, \hat{\phi}_i^-, c_i^+, c_i^-$ are stored for each player, its space complexity is in $\mathcal{O}(n)$. Moreover, SVARM is incremental and can be stopped at any time to return its estimates after executing line 20, or it can be run further with increased budget.

Theoretical Analysis. In the following we present theoretical results for SVARM. All proofs are given in Section C of the technical appendix. For the remainder of this section we assume that a minimum budget of $T \geq 2n + 2$ is given. This assumption guarantees the completion of the warm-up phase such that each positive and negative Shapley value estimate has at least one sample and an additional pair sampled in the loop. The lower bound on T is essentially twice the number of players n , which is a fairly weak assumption. We denote by $\bar{T} := T - 2n$ the number of time steps (budget) left after the warm-up phase. Moreover, we assume \bar{T} to be even for sake of simplicity such that a lower bound on the number of sampled pairs in the main part can be expressed

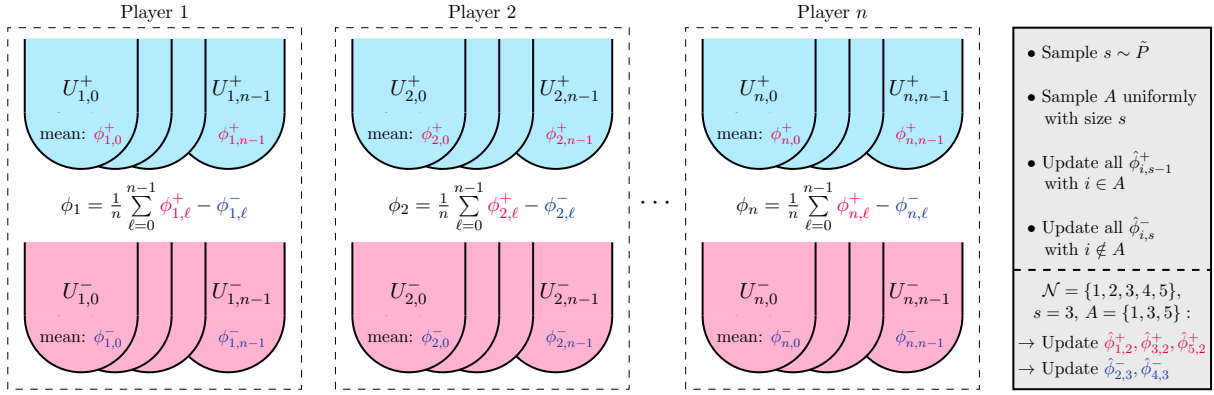


Figure 2: Illustration of Stratified SVARM’s sampling process and update rule: Each player i has urns $U_{i,\ell}^+ := \{S \cup \{i\} \mid S \subseteq \mathcal{N}_i, |S| = \ell\}$ and $U_{i,\ell}^- := \{S \mid S \subseteq \mathcal{N}_i, |S| = \ell\}$ for all $\ell \in \{0, \dots, n-1\}$, $2n$ in total, containing marbles which represent coalitions, with mean coalition worth $\phi_{i,\ell}^+$ and $\phi_{i,\ell}^-$. Stratified SVARM samples in each time step t a coalition $A_t \subseteq \mathcal{N}$ and updates the estimates of all players’ urns that contain the corresponding marble. Since each player’s urns form a partition of the powerset $\mathcal{P}(\mathcal{N})$, all players have exactly one urn updated with each sample.

by $\frac{T}{2} - n$. We begin with the unbiasedness of the estimates maintained by SVARM allowing us later to reduce the mean squared error (MSE) of each estimate to its variance.

Theorem 1. *The Shapley value estimate $\hat{\phi}_i$ of any $i \in \mathcal{N}$ obtained by SVARM is unbiased, i.e.,*

$$\mathbb{E}[\hat{\phi}_i] = \phi_i.$$

Next, we give a bound on the variance of each Shapley value estimate. For this purpose, we introduce notation for the variances of coalition values contained in ϕ_i^+ and ϕ_i^- . For a random set $A_i \subseteq \mathcal{N}_i$ distributed according to P^w let

$$\sigma_i^{+2} := \mathbb{V}[\nu(A_i \cup \{i\})] \text{ and } \sigma_i^{-2} := \mathbb{V}[\nu(A_i)]. \quad (8)$$

Theorem 2. *The variance of any player’s Shapley value estimate $\hat{\phi}_i$ obtained by SVARM is bounded by*

$$\mathbb{V}[\hat{\phi}_i] \leq \frac{2H_n}{T} (\sigma_i^{+2} + \sigma_i^{-2}).$$

Combining the unbiasedness in Theorem 1 with the latter variance bound implies the following result on the MSE.

Corollary 1. *The MSE of any player’s Shapley value estimate $\hat{\phi}_i$ obtained by SVARM is bounded by*

$$\mathbb{E}[(\hat{\phi}_i - \phi_i)^2] \leq \frac{2H_n}{T} (\sigma_i^{+2} + \sigma_i^{-2}).$$

Assuming that each variance term σ_i^{+2} and σ_i^{-2} is bounded by some constant independent of n (and T), the MSE bound in Corollary 1 is in $\mathcal{O}(\frac{\log n}{T-n})$ and so is the variance bound in Theorem 2. Note that this assumption is rather mild and satisfied if the underlying value function is bounded by constants independent of n , which again is the case for a wide range of games and in particular in explainable AI for global and local feature importance based on classification probabilities lying between 0 and 1. Further, as T is growing linearly with n by assumption, the denominator is essentially

driven by the asymptotics of T . Thus, the dependency on n is logarithmic, which is a significant improvement over existing theoretical results having a linear dependency on n like $\mathcal{O}(\frac{n}{T})$ for *ApproShapley* (Castro, Gómez, and Tejada 2009) or possibly worse (Simon and Vincent 2020). Finally, we present two probabilistic bounds on the approximated Shapley value. The first utilizes the variance bound shown in Theorem 2 by applying Chebyshev’s inequality.

Theorem 3. *The probability that the Shapley value estimate $\hat{\phi}_i$ of any fixed player $i \in \mathcal{N}$ deviates from ϕ_i by a margin of any fixed $\varepsilon > 0$ or greater is bounded by*

$$\mathbb{P}(|\hat{\phi}_i - \phi_i| \geq \varepsilon) \leq \frac{2H_n}{\varepsilon^2 T} (\sigma_i^{-2} + \sigma_i^{+2}).$$

The presented bound is in $\mathcal{O}(\frac{\log n}{T-n})$ and improves upon the bound derived by Chebyshev’s inequality of $\mathcal{O}(\frac{n}{T})$ for *ApproShapley* (Maleki et al. 2013). Our second bound derived by Hoeffding’s inequality is tighter, but requires the introduction of notation for the ranges of $\nu(A_i)$ and $\nu(A_i \cup \{i\})$:

$$r_i^+ := \max_{S \subseteq \mathcal{N}_i} \nu(S \cup \{i\}) - \min_{S \subseteq \mathcal{N}_i} \nu(S \cup \{i\}), \quad (9)$$

$$r_i^- := \max_{S \subseteq \mathcal{N}_i} \nu(S) - \min_{S \subseteq \mathcal{N}_i} \nu(S). \quad (10)$$

Theorem 4. *The probability that the Shapley value estimate $\hat{\phi}_i$ of any fixed player $i \in \mathcal{N}$ deviates from ϕ_i by a margin of any fixed $\varepsilon > 0$ or greater is bounded by*

$$\mathbb{P}(|\hat{\phi}_i - \phi_i| \geq \varepsilon) \leq 2e^{-\frac{\varepsilon^2}{4H_n^2}} + 4 \frac{e^{-\Psi \lfloor \frac{\varepsilon}{4H_n} \rfloor}}{e^\Psi - 1},$$

where $\Psi = 2\varepsilon^2 / (r_i^+ + r_i^-)^2$.

Note that this bound is exponentially decreasing with T and can be expressed asymptotically as $\mathcal{O}(e^{-\frac{\varepsilon^2}{(\log n)^2}})$. In comparison, the bounds of $\mathcal{O}(e^{-\frac{\varepsilon^2}{n}})$ for *ApproShapley*, $\mathcal{O}(ne^{-\frac{\varepsilon^2}{n^3}})$ for *Stratified Sampling* (Maleki et al. 2013), and the projected SGD variant (Simon and Vincent 2020) show worse asymptotic dependencies on n in comparison.

Stratified SVARM

On the basis of the representation of the Shapley value in Equation (4), we develop another approximation algorithm named *Stratified SVARM* to further pursue and reach the maximum sample reuse principle. Its crux is a refinement of SVARM obtained by stratifying the positive and the negative Shapley value ϕ_i^+ and ϕ_i^- . We exploit the latter to develop an even more powerful update rule that allows for updating all players simultaneously with each single coalition sampled. Both, ϕ_i^+ and ϕ_i^- can be rewritten using stratification such that each becomes an average of strata, whereas the strata themselves are averages of the coalitions' worth:

$$\phi_i^+ = \frac{1}{n} \sum_{\ell=0}^{n-1} \frac{1}{\binom{n-1}{\ell}} \sum_{\substack{S \subseteq \mathcal{N}_i \\ |S|=\ell}} \nu(S \cup \{i\}) =: \frac{1}{n} \sum_{\ell=0}^{n-1} \phi_{i,\ell}^+, \quad (11)$$

$$\phi_i^- = \frac{1}{n} \sum_{\ell=0}^{n-1} \frac{1}{\binom{n-1}{\ell}} \sum_{\substack{S \subseteq \mathcal{N}_i \\ |S|=\ell}} \nu(S) =: \frac{1}{n} \sum_{\ell=0}^{n-1} \phi_{i,\ell}^-. \quad (12)$$

We call $\phi_{i,\ell}^+$ the ℓ -th positive Shapley subvalue and $\phi_{i,\ell}^-$ the ℓ -th negative Shapley subvalue for all $\ell \in \mathcal{L} := \{0, \dots, n-1\}$. Now, we can write ϕ_i as

$$\phi_i = \frac{1}{n} \sum_{\ell=0}^{n-1} \phi_{i,\ell}^+ - \phi_{i,\ell}^-. \quad (13)$$

Note that this representation of ϕ_i coincides with Equation 6 in (Ancona, Öztireli, and Gross 2019). Intuitively speaking at the example of ϕ_i^+ (and analogously for ϕ_i^-), we partition the population of coalitions contained in ϕ_i^+ into n strata. Each *stratum* $\phi_{i,\ell}^+$ comprises all coalitions which include the player i and have cardinality $\ell + 1$. Instead of sampling directly for ϕ_i^+ , the stratification allows one to sample coalitions from each stratum, obtain mean estimates $\hat{\phi}_{i,\ell}^+$, and aggregate them to

$$\hat{\phi}_i^+ = \frac{1}{n} \sum_{\ell=0}^{n-1} \hat{\phi}_{i,\ell}^+ \quad (14)$$

in order to obtain an estimate for ϕ_i^+ . Due to the increase in homogeneity of the strata in comparison to their origin population, caused by the shared size and inclusion or exclusion of i for coalitions in the same stratum, one would expect the strata to have significantly lower variances and ranges resulting in approximations of better quality compared to SVARM. In combination with our bounds shown in Theorem 2 and Theorem 4, this should result in approximations of better quality. In the following we present further techniques for improvement which we apply for Stratified SVARM (Algorithm 2).

Exact Calculation. First, we observe that some strata contain very few coalitions. Thus, we calculate $\phi_{i,0}^+, \phi_{i,n-2}^+, \phi_{i,n-1}^+, \phi_{i,1}^-$, and $\phi_{i,n-1}^-$ for all players exactly by evaluating ν for all coalitions of size 1, $n-1$, and n . This requires $2n+1$ many evaluations of ν (see Algorithm 5 in Appendix B). We already obtain $\phi_{i,0}^- = \nu(\emptyset) = 0$

Algorithm 2: Stratified SVARM

Input: $\mathcal{N}, T \in \mathbb{N}$

- 1: $\hat{\phi}_{i,\ell}^+, \hat{\phi}_{i,\ell}^- \leftarrow 0$ for all $i \in \mathcal{N}$ and $\ell \in \mathcal{L}$
- 2: $c_{i,\ell}^+, c_{i,\ell}^- \leftarrow 0$ for all $i \in \mathcal{N}$ and $\ell \in \mathcal{L}$
- 3: EXACTCALCULATION(\mathcal{N})
- 4: WARMUP⁺(\mathcal{N})
- 5: WARMUP⁻(\mathcal{N})
- 6: $t \leftarrow 2n + 1 + 2 \sum_{s=2}^{n-2} \lceil \frac{n}{s} \rceil$
- 7: **while** $t < T$ **do**
- 8: Draw $s_t \sim \tilde{P}$
- 9: Draw A_t from $\{S \subseteq \mathcal{N} \mid |S| = s_t\}$ uniformly
- 10: UPDATE(A_t)
- 11: $t \leftarrow t + 1$
- 12: **end while**
- 13: $\hat{\phi}_i \leftarrow \frac{1}{n} \sum_{\ell=0}^{n-1} \hat{\phi}_{i,\ell}^+ - \hat{\phi}_{i,\ell}^-$ for all $i \in \mathcal{N}$

Output: $\hat{\phi}_1, \dots, \hat{\phi}_n$

by definition. As a consequence, we can exclude the sizes 0, 1, $n-1$, and n from further consideration. We assume for the remainder that $n \geq 4$, otherwise we would have already calculated all Shapley values exactly.

Refined Warm-Up. Next, we split the warm-up into two parts, one for the positive, the other for the negative Shapley subvalues (see Algorithm 6 and 7 in Appendix B). Each collects for each estimate $\hat{\phi}_{i,\ell}^+$ or $\hat{\phi}_{i,\ell}^-$, respectively, one sample and consumes a budget of $\sum_{s=2}^{n-2} \lceil \frac{n}{s} \rceil$.

Enhanced Update Rule. Thanks to the stratified representation of the Shapley value, we can enhance SVARM's update rule and update with each sampled coalition $A_t \subseteq \mathcal{N}$ the estimates $\hat{\phi}_{i,|A_t|-1}^+$ for all $i \in A_t$ and $\hat{\phi}_{i,|A_t|}^-$ for all $i \notin A_t$. Thus, we can update all estimates $\hat{\phi}_i$ at once with a single sample. This enhanced update step is given in Algorithm 4 (see Appendix B) and illustrated in Figure 2. In order to obtain unbiased estimates, it suffices to select an arbitrary size s of the coalition A to be sampled and draw A uniformly at random from the set of coalitions with size s . We go one step further and choose not only the coalition A , but also the size s randomly according to a specifically tailored probability distribution \tilde{P} over $\{2, \dots, n-2\}$, which leads to simpler bounds in our theoretical analysis in which each stratum receives the same weight. We define for n even:

$$\tilde{P}(s) := \begin{cases} \frac{n \log n - 1}{2sn \log n \left(H_{\frac{n}{2}-1} - 1 \right)} & \text{if } s \leq \frac{n-2}{2} \\ \frac{1}{n \log n} & \text{if } s = \frac{n}{2} \\ \frac{n \log n - 1}{2(n-s)n \log n \left(H_{\frac{n}{2}-1} - 1 \right)} & \text{otherwise} \end{cases},$$

$$\text{and for } n \text{ odd: } \tilde{P}(s) := \begin{cases} \frac{1}{2s \left(H_{\frac{n-1}{2}-1} - 1 \right)} & \text{if } s \leq \frac{n-1}{2} \\ \frac{1}{2(n-s) \left(H_{\frac{n-1}{2}-1} - 1 \right)} & \text{otherwise} \end{cases}.$$

Note that Stratified SVARM is incremental just as SVARM, but in contrast, requires quadratic space $\mathcal{O}(n^2)$ as it stores estimates and counters for each player *and* stratum.

Theoretical Analysis. Similar to SVARM, we present in the following our theoretical results for Stratified SVARM. All proofs are given in Appendix D. Again, we assume a minimum budget of $T \geq 2n + 1 + 2 \sum_{s=2}^{n-2} \lceil \frac{n}{s} \rceil =: W \in \mathcal{O}(n \log n)$, guaranteeing the completion of the warm-up phase, and denote by $\bar{T} = T - W$ the budget left after the warm-up phase. We start by showing that Stratified SVARM is not afflicted with any bias.

Theorem 5. *The Shapley value estimate $\hat{\phi}_i$ of any $i \in \mathcal{N}$ obtained by Stratified SVARM is unbiased, i.e.,*

$$\mathbb{E}[\hat{\phi}_i] = \phi_i.$$

Next, we consider the variance of the Shapley value estimates and quickly introduce some notation. Let $A_{i,\ell} \subseteq \mathcal{N}_i$ be a random coalition of size ℓ distributed with $\mathbb{P}(A_{i,\ell} = S) = \binom{n-1}{\ell}^{-1}$. Define the strata variances

$$\sigma_{i,\ell}^{+2} := \mathbb{V}[\nu(A_{i,\ell} \cup \{i\})] \text{ and } \sigma_{i,\ell}^{-2} := \mathbb{V}[\nu(A_{i,\ell})]. \quad (15)$$

Theorem 6. *The variance of any player's Shapley value estimate $\hat{\phi}_i$ obtained by Stratified SVARM is bounded by*

$$\mathbb{V}[\hat{\phi}_i] \leq \frac{2 \log n}{n\bar{T}} \sum_{\ell=1}^{n-3} \sigma_{i,\ell}^{+2} + \sigma_{i,\ell+1}^{-2}.$$

Together with the unbiasedness shown in Theorem 5, the variance bound implies the following MSE bound.

Corollary 2. *The MSE of any player's Shapley value estimate $\hat{\phi}_i$ obtained by Stratified SVARM is bounded by*

$$\mathbb{E}[(\hat{\phi}_i - \phi_i)^2] \leq \frac{2 \log n}{n\bar{T}} \sum_{\ell=1}^{n-3} \sigma_{i,\ell}^{+2} + \sigma_{i,\ell+1}^{-2}.$$

With our choice of the sampling distribution \tilde{P} we achieved an easily interpretable bound on the MSE in which each stratum variance is equally weighted. Assuming that each stratum variance is bounded by some constant independent of n , the MSE bound in Corollary 2 is in $\mathcal{O}(\frac{\log n}{T-n \log n})$. Note that, by assumption, T is growing log-linearly with n so that the denominator is essentially driven by the asymptotics of T . Again, compared to existing theoretical results, with linear dependence on n , the logarithmic dependence on n is a significant improvement. Still, it is worth emphasizing that the more homogeneous strata with lower variances constitute the core improvement of Stratified SVARM, which are not reflected within the \mathcal{O} -notation. Our first probabilistic bound is obtained by Chebyshev's inequality and the bound from Theorem 6.

Theorem 7. *The probability that the Shapley value estimate $\hat{\phi}_i$ of any fixed player $i \in \mathcal{N}$ deviates from ϕ_i by a margin of any fixed $\varepsilon > 0$ or greater is bounded by*

$$\mathbb{P}(|\hat{\phi}_i - \phi_i| \geq \varepsilon) \leq \frac{2 \log n}{\varepsilon^2 n \bar{T}} \sum_{\ell=1}^{n-3} \sigma_{i,\ell}^{+2} + \sigma_{i,\ell+1}^{-2}.$$

Lastly, our second probabilistic bound derived via Hoeffding's inequality is tighter, but less trivial. It requires some further notation, namely the ranges of the strata values:

$$r_{i,\ell}^+ := \max_{S \subseteq \mathcal{N}_i: |S|=\ell} \nu(S \cup \{i\}) - \min_{S \subseteq \mathcal{N}_i: |S|=\ell} \nu(S \cup \{i\}), \quad (16)$$

$$r_{i,\ell}^- := \max_{S \subseteq \mathcal{N}_i: |S|=\ell} \nu(S) - \min_{S \subseteq \mathcal{N}_i: |S|=\ell} \nu(S). \quad (17)$$

Theorem 8. *The probability that the Shapley value estimate $\hat{\phi}_i$ of any fixed player $i \in \mathcal{N}$ deviates from ϕ_i by a margin of any fixed $\varepsilon > 0$ or greater is bounded by $\mathbb{P}(|\hat{\phi}_i - \phi_i| \geq \varepsilon)$*

$$\leq 2(n-3) \left(e^{-\frac{\bar{T}}{8n^2(\log n)^2}} + 2 \frac{e^{-\Psi \lfloor \frac{\bar{T}}{4n \log n} \rfloor}}{e^\Psi - 1} \right),$$

$$\text{where } \Psi = 2\varepsilon^2 n^2 / \left(\sum_{\ell=1}^{n-3} r_{i,\ell}^+ + r_{i,\ell+1}^- \right)^2.$$

This bound is of order $\mathcal{O}(ne^{-\frac{T-n \log n}{n^2(\log n)^2}})$ showing a slightly worse dependency on n compared to Theorem 4 due to the introduction of strata.

Empirical Results

To complement our theoretical findings, we evaluate our algorithms and its competitors on commonly considered synthetic cooperative games and explainable AI scenarios in which Shapley values need to be approximated. In particular, we select parameterless algorithms that do not rely on provided knowledge about the value function of the problem instance at hand, since ours do not either. Besides the sampling distribution \tilde{P} over coalition sizes proposed for Stratified SVARM (S-SVARM), we also consider sampling with the simpler uniform distribution over all sizes from 2 to $n-2$ (S-SVARM uniform). In order to allow for a fair comparison with KernelSHAP, which samples coalitions without replacement, we include with S-SVARM⁺ (uniform) an empirical version of S-SVARM without the warm-up that also samples without replacement to compensate for this underlying advantage (see Algorithm 8 in Appendix B), which obviously comes at the price of space complexity linear in T .

We run the algorithms multiple times on the selected game types and measure their performances by the mean squared error (MSE) averaged over all players and runs depending on a range of fixed budget values T . Measuring the approximation quality by the MSE requires the true Shapley values of the considered games to be available. These are either given by a polynomial closed-form solution for the synthetic games (see Section 6.1) or we compute them exhaustively for our explanation tasks (see Section 6.2). The results of our evaluation are shown in Figure 3 and are presented in more detail in Appendix F.

As already said, we judge the algorithms' approximation qualities in dependence on the spent budget (model evaluations) T instead of the consumed runtime. In fact, the algorithms differ in actual runtime. For example SVARM performs less arithmetic operations than Stratified SVARM since it does not update all players' estimates $\hat{\phi}_i^+$ or $\hat{\phi}_i^-$

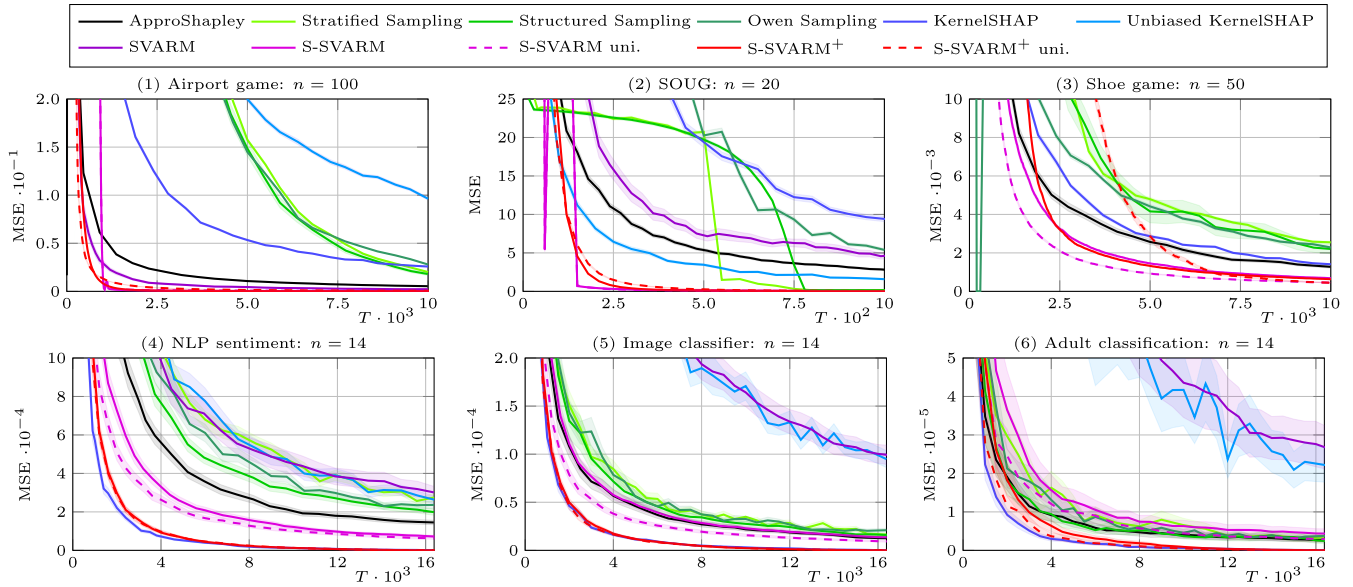


Figure 3: Averaged MSE and standard errors over 100 repetitions in dependence of fixed budget T : (1) Airport game, (2) Shoe game, (3) SOUG game, (4) NLP sentiment analysis, (5) Image classifier, (6) Adult classification.

with each sample. Some algorithms, e.g. KernelSHAP, vary strongly in their time consumption per sample since a costly quadratic optimization problem needs to be solved after observing all samples. We intentionally avoid the runtime comparison for three reasons: (i) the observed runtimes may differ depending on the actual implementation, (ii) the fixed-budget setting facilitates a coherent theoretical analysis where the observed information is restricted, (iii) evaluating the worth of a coalition poses the bottleneck in explanation tasks, rendering the difference in performed arithmetic operations negligible.

Synthetic Games

Cooperative games with polynomial closed-form solutions of their Shapley values are well suited for tracking the approximation error for large player numbers. We exploit this fact and investigate a broad range of player numbers n which are significantly higher than those for the explanation tasks. We conduct experiments on the predefined Shoe and Airport game as done in (Castro, Gómez, and Tejada 2009; Castro et al. 2017). Their degree of non-additivity poses a difficult challenge to all approximation algorithms. Further, we consider randomly generated Sum of Unanimity Games (SOUG) games (van Campen et al. 2018) which are capable of representing any cooperative game. The value function and Shapley values of each game are given in Appendix E.

We observe that S-SVARM itself already shows reliably good approximation performance across all considered games and budget ranges. It is significantly superior to its competitors ApproShapley and KernelSHAP and as expected, S-SVARM⁺ extends the lead in approximation quality even more. In contrast, SVARM can rarely keep up with its refined counterpart S-SVARM. However, in light of the bounds on the MSEs in Corollary 1 and 2 this is not surpris-

ing: SVARM’s MSE bound scales linearly with the variances σ_i^{+2} and σ_i^{-2} of all coalition values containing respectively not containing i , while the relevant variance terms $\sigma_{i,\ell}^{+2}$ and $\sigma_{i,\ell}^{-2}$ for S-SVARM are restricted to coalitions of fixed size. In most games, the latter terms are significantly lower since coalitions of the same size are plausibly closer in worth. Finally, S-SVARM is quite robust regarding the magnitude of the standard errors.

Explainability Games

We further conduct experiments on cooperative games stemming from real-world explainability scenarios, in particular, use cases in which local feature importance of machine learning models are to be quantified via Shapley values. The NLP sentiment analysis game is based on the DistilBERT (Sanh et al. 2019) model architecture and consists of randomly selected movie reviews from the IMDB dataset (Maas et al. 2011) containing 14 words. Missing features are masked in the tokenized representation and the value of a set is its sentiment score. In the image classifier game, we explain the output of a ResNet18 (He et al. 2016) trained on ImageNet (Deng et al. 2009). The images’ pixels are summarized into $n = 14$ super-pixels and absent features are masked with mean imputation. The worth of a coalition is the returned class probability of the model (using only the present super-pixels) for the class of the original prediction which was made with all pixels being present. For the adult classification game, we train a gradient-boosted tree model on the adult dataset (Becker and Kohavi 1996). A coalition’s worth is the predicted class probability of the true income class (income above or below 50 000) of the given datapoint with the absent features being removed via mean imputation. Since no polynomial closed-form solution exists for

the Shapley values in these games, we compute them exhaustively, limiting us to a feasible number of players for which we can track the MSE. While this restricts us to a player number (tokens, superpixels, features) of $n = 14$ due to limited computational resources, this is arguably still an appropriate and commonly appearing number of entities involved in an explanation task. We refer to Appendix E for a more detailed explanation of the chosen games.

A first observation is the close head-to-head race between S-SVARM⁺ and KernelSHAP across the considered games leaving all other methods behind. Thus, S-SVARM⁺ is the first sample-mean-based approach achieving rivaling state-of-the-art approximation quality. KernelSHAP’s counterpart Unbiased KernelSHAP, designed to facilitate approximation guarantees similar to our theoretical results which KernelSHAP lacks, is clearly outperformed by S-SVARM. Given the consistency demonstrated by S-SVARM and S-SVARM⁺, we claim that both constitute a reliable choice under absence of domain knowledge. We conjecture that the reason for the slight performance decrease of S-SVARM from synthetic to explainability games lies not only within the latent structure of ν , but is also caused by the lower player numbers. As our theoretical results indicate, its sample efficiency grows with n due to its enhanced update rule. However, conducting experiments with larger n becomes computationally prohibitive for explainability games, since the Shapley values have to be calculated exhaustively in order to track the approximation error. Further, our results indicate the robustness of S-SVARM(+) w.r.t. the utilized distribution \tilde{P} , which allows us to use the uniform distribution without performance loss, and secondly shows that our derived distribution is not just a theoretical artifact, but a valid contribution to express simpler bounds which are easier to grasp and interpret.

Conclusion

We considered the problem of precisely approximating the Shapley value of all players in a cooperative game under the restriction that the value function can be evaluated only a given number of times. We presented a reformulation of the Shapley value, detached from the ubiquitous notion of marginal contribution, facilitating the approximation by estimates of which a multitude can be updated with each access to the value function. On this basis, we proposed two approximation algorithms, SVARM and Stratified SVARM, which have a number of desirable properties. Both are parameter-free, incremental, domain-independent, unbiased, and do not require any prior knowledge of the value function. Further, Stratified SVARM shows a satisfying compromise between peak approximation quality and consistency across all considered games, paired with unmatched theoretical guarantees regarding its approximation quality. While fulfilling more desirable properties and not having to solve a quadratic optimization problem of size T in comparison to the state-of-the-art method KernelSHAP, effectively disabling on-the-fly approximations, our simpler sample-mean-based method Stratified SVARM⁺ can fully keep up in common explainable AI scenarios, and even shows empirical su-

periority on synthetic games.

Limitations and Future Work. The quadratically growing number of strata w.r.t. n might pose a challenge for higher player numbers, which future work could remedy by applying a coarser stratification that assigns multiple coalition sizes to a single stratum. One could investigate the empirical behavior in further popular explanation domains such as data valuation, federated learning, or neuron importance and extend our evaluation to scenarios with higher player numbers. Since the true Shapley values are not accessible for larger n , a different measure of approximation quality than the MSE needs to be taken for reference. The convergence speed of the estimates is a naturally arising alternative. Our empirical results give further evidence for the non-existence of a universally best approximation algorithm and encourage future research into the cause of the observed differences in performance w.r.t. the game type. Further, it would be interesting to analyze whether structural properties of the value function, such as monotonicity or submodularity, have an impact on the approximation quality of both algorithms.

Acknowledgments

This research was supported by the research training group Dataninja (Trustworthy AI for Seamless Problem Solving: Next Generation Intelligence Joins Robust Data Analysis) funded by the German federal state of North Rhine-Westphalia. We gratefully acknowledge funding by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation): TRR 318/1 2021 – 438445824. We would like to thank Fabian Fumagalli and especially Patrick Becker for their efforts in supporting our implementation.

References

- Ancona, M.; Öztireli, C.; and Gross, M. H. 2019. Explaining Deep Neural Networks with a Polynomial Time Algorithm for Shapley Value Approximation. In *Proceedings of the 36th International Conference on Machine Learning*, volume 97, 272–281.
- Becker, B.; and Kohavi, R. 1996. Adult. UCI Machine Learning Repository.
- Burgess, M. A.; and Chapman, A. C. 2021. Approximating the Shapley Value Using Stratified Empirical Bernstein Sampling. In *Proceedings of the Thirtieth International Joint Conference on Artificial Intelligence*, 73–81.
- Castro, J.; Gómez, D.; Molina, E.; and Tejada, J. 2017. Improving Polynomial Estimation of the Shapley Value by Stratified Random Sampling with Optimum Allocation. *Computers & Operations Research*, 82: 180–188.
- Castro, J.; Gómez, D.; and Tejada, J. 2009. Polynomial Calculation of the Shapley Value based on Sampling. *Computers & Operations Research*, 36(5): 1726–1730.
- Chen, H.; Covert, I. C.; Lundberg, S. M.; and Lee, S.-I. 2023. Algorithms to Estimate Shapley Value Feature Attributions. *Nature Machine Intelligence*, 5: 590–601.
- Cohen, S. B.; Dror, G.; and Ruppin, E. 2007. Feature Selection via Coalitional Game Theory. *Neural Computation*, 19(7): 1939–1961.

- Corder, K.; and Decker, K. 2019. Shapley Value Approximation with Divisive Clustering. In *18th IEEE International Conference On Machine Learning And Applications*, 234–239.
- Covert, I.; and Lee, S.-I. 2021. Improving KernelSHAP: Practical Shapley Value Estimation Using Linear Regression. In *24th International Conference on Artificial Intelligence and Statistics*, volume 130 of *PMLR*, 3457–3465.
- Covert, I.; Lundberg, S.; and Lee, S.-I. 2019. Shapley Feature Utility. In *Machine Learning in Computational Biology*.
- Deng, J.; Dong, W.; Socher, R.; Li, L.; Li, K.; and Fei-Fei, L. 2009. ImageNet: A Large-Scale Hierarchical Image Database. In *IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, 248–255.
- Deng, X.; and Papadimitriou, C. H. 1994. On the Complexity of Cooperative Solution Concepts. *Mathematics of Operations Research*, 19(2): 257–266.
- Ghorbani, A.; and Zou, J. Y. 2019. Data Shapley: Equitable Valuation of Data for Machine Learning. In *Proceedings of the 36th International Conference on Machine Learning*, volume 97, 2242–2251.
- Ghorbani, A.; and Zou, J. Y. 2020. Neuron Shapley: Discovering the Responsible Neurons. In Larochelle, H.; Ranzato, M.; Hadsell, R.; Balcan, M.; and Lin, H., eds., *Advances in Neural Information Processing Systems*, volume 33.
- He, K.; Zhang, X.; Ren, S.; and Sun, J. 2016. Deep Residual Learning for Image Recognition. In *IEEE Conference on Computer Vision and Pattern Recognition*, 770–778.
- Illés, F.; and Kerényi, P. 2019. Estimation of the Shapley Value by Ergodic Sampling. *CoRR*, abs/1906.05224.
- Lomeli, M.; Rowland, M.; Gretton, A.; and Ghahramani, Z. 2019. Antithetic and Monte Carlo Kernel Estimators for Partial Rankings. *Statistics and Computing*, 29(5): 1127–1147.
- Lundberg, S. M.; and Lee, S.-I. 2017. A Unified Approach to Interpreting Model Predictions. In *Advances in Neural Information Processing Systems*, volume 30, 4768–4777.
- Luo, C.; Zhou, X.; and Lev, B. 2022. Core, Shapley Value, Nucleolus and Nash Bargaining Solution: A Survey of Recent Developments and Applications in Operations Management. *Omega*, 110: 102638.
- Maas, A. L.; Daly, R. E.; Pham, P. T.; Huang, D.; Ng, A. Y.; and Potts, C. 2011. Learning Word Vectors for Sentiment Analysis. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, 142–150.
- Maleki, S.; Tran-Thanh, L.; Hines, G.; Rahwan, T.; and Rogers, A. 2013. Bounding the Estimation Error of Sampling-based Shapley Value Approximation With/Without Stratifying. *CoRR*, abs/1306.4265.
- Mitchell, R.; Cooper, J.; Frank, E.; and Holmes, G. 2022. Sampling Permutations for Shapley Value Estimation. *Journal of Machine Learning Research*, 23(43): 1–46.
- O’Brien, G.; Gamal, A. E.; and Rajagopal, R. 2015. Shapley Value Estimation for Compensation of Participants in Demand Response Programs. *IEEE Transactions on Smart Grid*, 6(6): 2837–2844.
- Okhrati, R.; and Lipani, A. 2020. A Multilinear Sampling Algorithm to Estimate Shapley Values. In *25th International Conference on Pattern Recognition*, 7992–7999.
- Owen, G. 1972. Multilinear Extensions of Games. *Management Science*, 18: 64–79.
- Rozemberczki, B.; and Sarkar, R. 2021. The Shapley Value of Classifiers in Ensemble Games. In *30th ACM International Conference on Information and Knowledge Management*, 1558–1567.
- Rozemberczki, B.; Watson, L.; Bayer, P.; Yang, H.-T.; Kiss, O.; Nilsson, S.; and Sarkar, R. 2022. The Shapley Value in Machine Learning. In *Proceedings of the 31st International Joint Conference on Artificial Intelligence*, 5572–5579.
- Rubinstein, R. Y.; and Kroese, D. P. 2016. *Simulation and the Monte Carlo Method*. John Wiley & Sons.
- Sanh, V.; Debut, L.; Chaumond, J.; and Wolf, T. 2019. DistilBERT, a Distilled Version of BERT: Smaller, Faster, Cheaper and Lighter. *CoRR*, abs/1910.01108.
- Shapley, L. S. 1953. A Value for n-Person Games. In *Contributions to the Theory of Games, Volume II*, 307–318. Princeton University Press.
- Simon, G.; and Vincent, T. 2020. A Projected Stochastic Gradient Algorithm for Estimating Shapley Value Applied in Attribute Importance. In *Machine Learning and Knowledge Extraction*, 97–115.
- Soufiani, H. A.; Chickering, D. M.; Charles, D. X.; and Parkes, D. C. 2014. Approximating the Shapley Value via Multi-Issue Decompositions. In *Proceedings of the International conference on Autonomous Agents and Multi-Agent Systems*, volume 2.
- Touati, S.; Radjef, M. S.; and Sais, L. 2021. A Bayesian Monte Carlo Method for Computing the Shapley Value: Application to Weighted Voting and Bin Packing Games. *Computers & Operations Research*, 125: 105094.
- van Campen, T.; Hamers, H.; Husslage, B.; and Lindelauf, R. 2018. A New Approximation Method for the Shapley Value Applied to the WTC 9/11 Terrorist Attack. *Social Network Analysis and Mining*, 8(3): 1–12.
- Wang, J. T.; and Jia, R. 2023. Data Banzhaf: A Robust Data Valuation Framework for Machine Learning. In *26th International Conference on Artificial Intelligence and Statistics*, volume 206 of *PMLR*, 6388–6421.