

# Approximating the Shapley Value without Marginal Contributions

Patrick Kolpaczki<sup>1</sup>, Viktor Bengs<sup>2</sup>,  
Maximilian Muschalik<sup>2</sup>, and Eyke Hüllermeier<sup>2</sup>

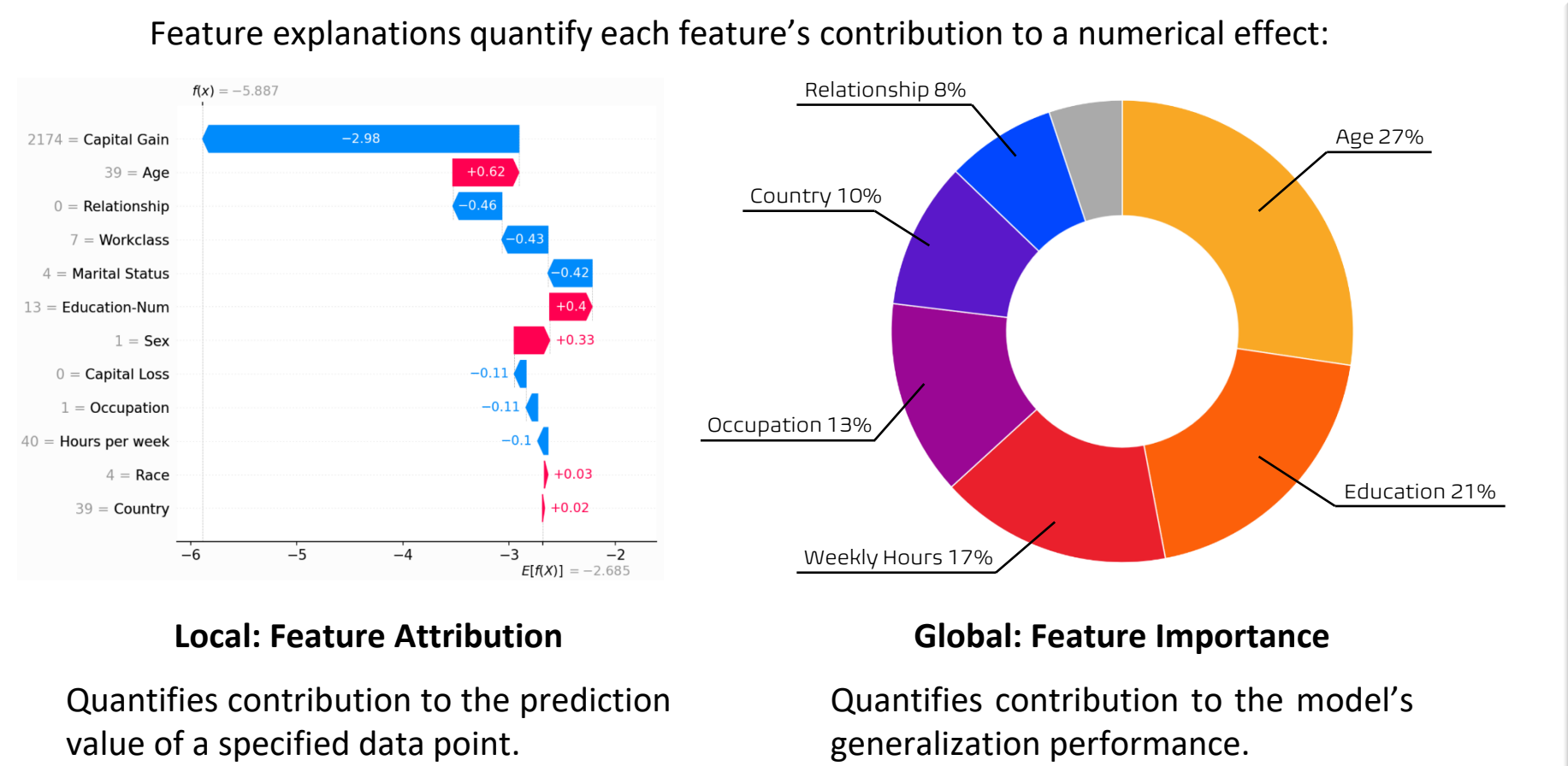


1) Paderborn University, Paderborn, Germany



2) LMU Munich, Munich, Germany

## Motivation: Additive Feature Explanations

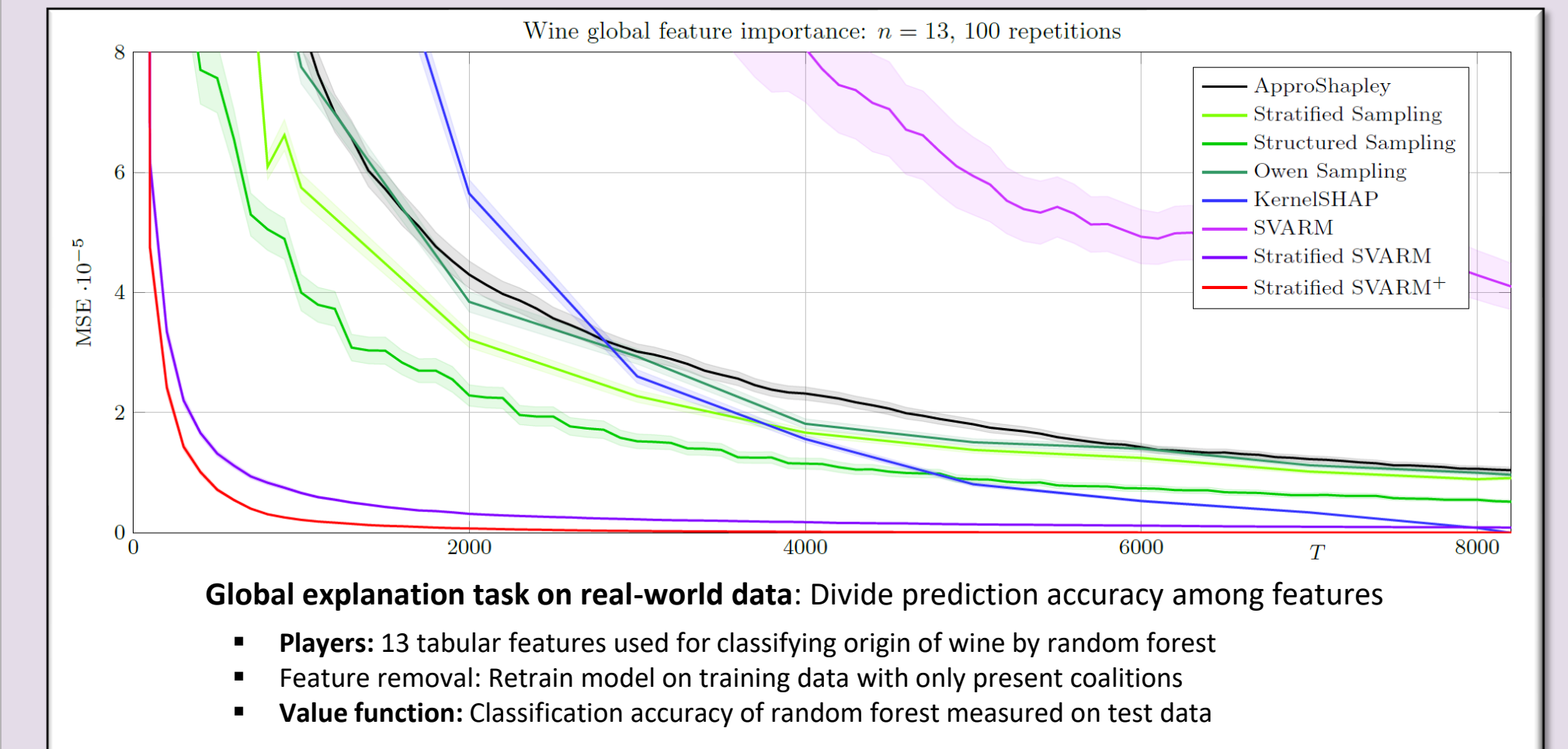


## Contribution

We propose a novel approximation algorithm for the Shapley value with:

- New combination of stratified representation + update mechanism
  - Novel strong theoretical guarantees
  - State-of-the-art empirical performance
- Desirable properties:**
- ✓ Model-agnostic / domain-independent  
→ Applicable for data valuation, neuron importance, etc. and even outside of ML
  - ✓ No hyperparameters  
→ No fine-tuning
  - ✓ Estimates available at any time  
→ Budget can be cut and extended arbitrarily
  - ✓ Uncertainty-aware  
→ Allows construction of confidence intervals

## Empirical Evaluation



## Shapley Value

- **Player set**  $\mathcal{N} = \{1, \dots, n\}$  → Features, datapoints, neurons, base learners etc.
- **Value function**  $v: \mathcal{P}(\mathcal{N}) \rightarrow \mathbb{R}$  → Predicted value, generalization performance with  $v(\emptyset) = 0$

**Definition: Shapley Value** (Shapley, 1953)

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

$\Delta_i(S)$

Marginal contribution  $\Delta_i(S)$ : Increase in collective benefit when  $i$  joins  $S$ .

- Unique solution to fulfill desirable axioms: Efficiency, Symmetry, Additivity, Null-Property
- Computational effort scales **exponentially** with  $n$ :  $2^n$  coalitions in total

### Fixed-budget approximation problem:

- Given cooperative game  $(N, v)$  with unknown Shapley values  $\phi_1, \dots, \phi_n$
- Budget  $T$ : Allowed number of evaluations of  $v$  (bottleneck due to model access)  
Model evaluations (inference, retraining) pose bottleneck on runtime rather than arithmetic operations
- Minimize mean squared error (MSE) averaged over all players for estimates  $\hat{\phi}_1, \dots, \hat{\phi}_n$ :

$$\frac{1}{n} \sum_{i=1}^n \mathbb{E} [(\hat{\phi}_i - \phi_i)^2]$$

### Approximation by sampling marginal contributions:

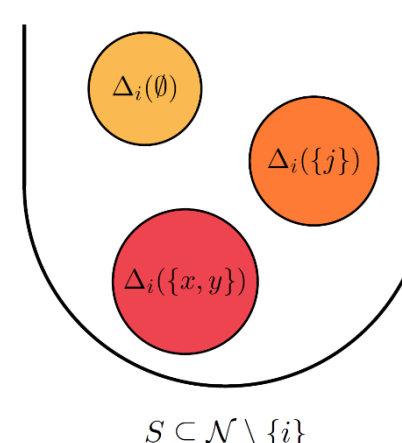
- Weights form a well-defined probability distribution:

$$\sum_{S \subseteq \mathcal{N} \setminus \{i\}} \frac{1}{n \cdot \binom{n-1}{|S|}} = 1$$

→ The Shapley value is the expected marginal contribution:

$$\phi_i = \mathbb{E} [\Delta_i(S)]$$

- Obtain  $\hat{\phi}_i$  by sampling marginal contributions according to weights



- One separate approximation problem for each player

**Problem:** Notion of marginal contributions is inefficient!

One update of  $\hat{\phi}_i$  with  $\Delta_i(S) = v(S \cup \{i\}) - v(S)$  costs **2 budget tokens**

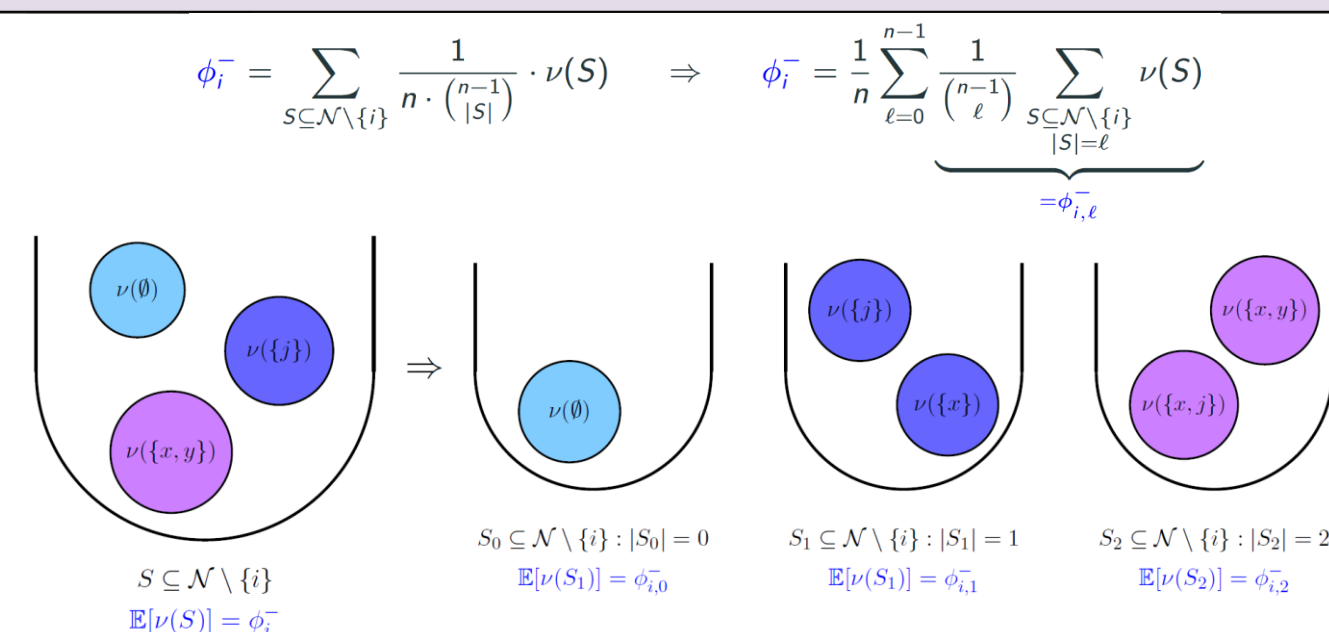
## Approximation Algorithm

### Stratified Representation

$$\phi_i = \underbrace{\frac{1}{n} \sum_{\ell=0}^{n-1} \frac{1}{\binom{n-1}{\ell}} \sum_{S \subseteq \mathcal{N} \setminus \{i\}} \nu(S \cup \{i\})}_{=\phi_{i,\ell}^+} - \underbrace{\frac{1}{n} \sum_{\ell=0}^{n-1} \frac{1}{\binom{n-1}{\ell}} \sum_{S \subseteq \mathcal{N} \setminus \{i\}} \nu(S)}_{=\phi_{i,\ell}^-} = \frac{1}{n} \sum_{\ell=0}^{n-1} \phi_{i,\ell}^+ - \phi_{i,\ell}^-$$

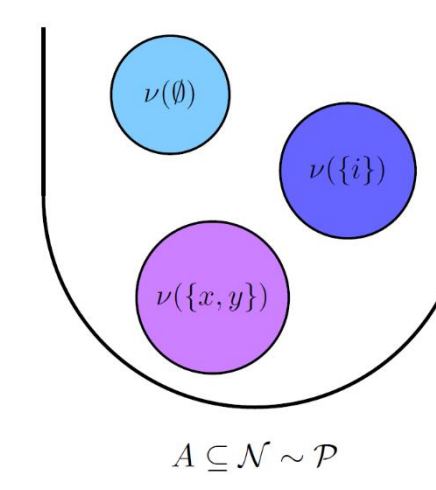
### Stratification of $\phi_i^+$ and $\phi_i^-$ :

- Partitions the marbles into  $n$  many strata
- Strata are grouped by size
- More homogeneous than base population
- Enables enhanced update mechanism
- Maintain estimates  $\hat{\phi}_{i,\ell}^+, \hat{\phi}_{i,\ell}^-$



### Stratified Shapley Value Approximation without Requesting Marginals (Stratified SVARM)

- Calculate all strata exactly with coalitions of size  $0, 1, \dots, n-1, n$
- Perform warmup to initialize all estimates:  $W \in O(n \log n)$
- Repeat with remaining budget  $\bar{T} = T - W$ :
  - Draw coalition size  $s \in \{2, \dots, n-2\} \sim \tilde{P}(s)$
  - Draw coalition of size  $s$  uniformly at random and evaluate  $v(A)$
  - Update  $\hat{\phi}_{i,s-1}^+$  for all  $i \in A$
  - Update  $\hat{\phi}_{i,s}^-$  for all  $i \notin A$



### Theorem 6. & Corollary 2. Variance and MSE

The variance and MSE of any estimate  $\hat{\phi}_i$  returned by Stratified SVARM is bounded by

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

with stratum variances  $\sigma_{i,\ell}^+ = \mathbb{V}[\nu(A \cup \{i\})]$  and  $\sigma_{i,\ell}^- = \mathbb{V}[\nu(A)]$  for  $A \subseteq \mathcal{N} \setminus \{i\}$  with  $|A| = \ell$  drawn u.a.r.

### Theorem 5. Unbiasedness

The estimate  $\hat{\phi}_i$  of any  $i \in \mathcal{N}$  returned by Stratified SVARM is unbiased, i.e.,  $\mathbb{E} [\hat{\phi}_i] = \phi_i, \forall i \in \mathcal{N}$ .

### Theorem 7. PAC bound

For any estimate  $\hat{\phi}_i$  returned by Stratified SVARM and  $\epsilon > 0$  holds

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

