

# Explaining probabilistic predictions on the simplex with Shapley compositions

**Paul-Gauthier Noé**<sup>1</sup>, Miquel Perelló-Nieto<sup>2</sup>,  
Peter Flach<sup>2</sup>, Jean-François Bonastre<sup>1</sup>

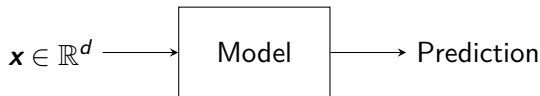
<sup>1</sup>Avignon Université

<sup>2</sup>University of Bristol

LIA seminar, Avignon Université  
January 18, 2024

# Introduction

## Local explanation in machine learning:



*Given one instance  $\mathbf{x}$  with  $d$  features, what is the contribution/effect of a feature's value on the prediction?*

$\neq$  Global explanation

# Introduction

Examples of local explanation methods:

- Local Interpretable Model-Agnostic Explanations (LIME)<sup>1</sup>,
- Shapley values<sup>2</sup> (SHAP toolkit<sup>3</sup>)

---

<sup>1</sup>Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. “” Why should i trust you?” Explaining the predictions of any classifier”. In: *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*. 2016, pp. 1135–1144.

<sup>2</sup>Erik Štrumbelj and Igor Kononenko. “Explaining prediction models and individual predictions with feature contributions”. In: *Knowledge and information systems* 41 (2014), pp. 647–665.

<sup>3</sup>Scott M Lundberg and Su-In Lee. “A Unified Approach to Interpreting Model Predictions”. In: *Advances in Neural Information Processing Systems* 30. Ed. by I. Guyon et al. Curran Associates, Inc., 2017, pp. 4765–4774.

# Introduction

## Shapley values in cooperative game theory<sup>4</sup>

- Distributes the total payoff among the players.
- The unique quantity respecting a set of desired axiomatic properties:
  - ▶ Linearity:

$$\phi_{\alpha v + (1-\alpha)w}(i) = \alpha \phi_v(i) + (1-\alpha) \phi_w(i), \quad (1)$$

for a player  $i$  and two games  $v$  and  $w$ , and for  $\alpha \in [0, 1]$


- ▶ Efficiency,

$$\sum_{i \in C} \phi_v(i) = v(C), \quad (2)$$

(the sum of the value is equal to the total payoff)

- ▶ Symmetry

---

<sup>4</sup>Lloyd S Shapley et al. "A value for n-person games". In: (1953), pp 307-317. 

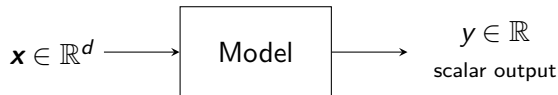
# Introduction

## Shapley values in machine learning

# Introduction

## Shapley values in machine learning

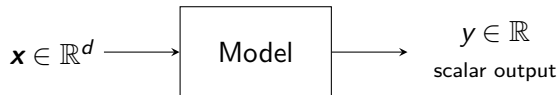
- Features are treated as players, and the scalar output of the model as the payoff,



# Introduction

## Shapley values in machine learning

- Features are treated as players, and the scalar output of the model as the payoff,

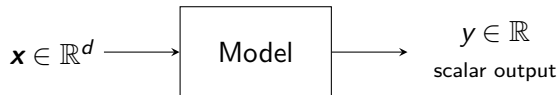


- Binary classifier, regressor with one-dimensional output ✓

# Introduction

## Shapley values in machine learning

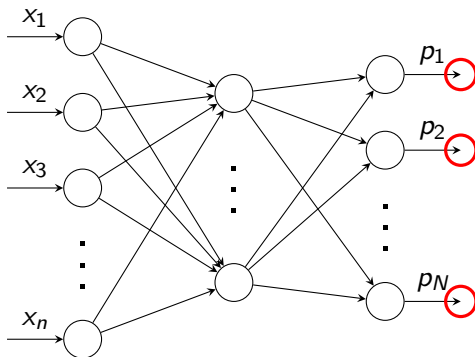
- Features are treated as players, and the scalar output of the model as the payoff,



- Binary classifier, regressor with one-dimensional output ✓
- Multiclass classifier ✗  
ex: The output of a softmax layer lives on a multidimensional simplex!

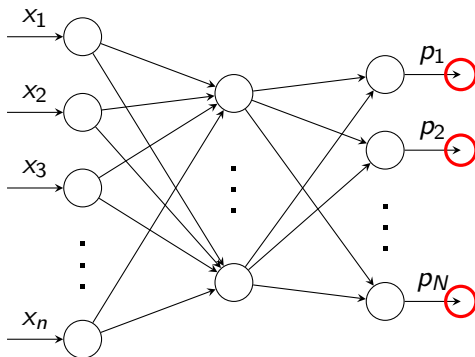


# Introduction



Some explain the output one-by-one,

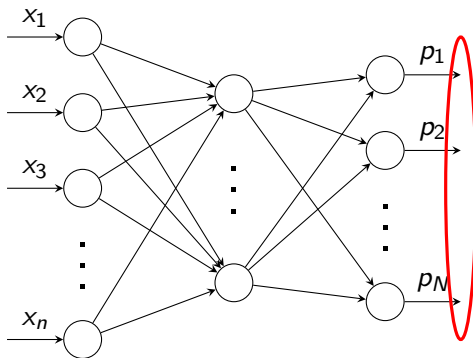
# Introduction



Some explain the output one-by-one,

But a probability distribution lives on a simplex,  
**The relative information matter!!**

# Introduction



We will explain the probabilities all together using the *Aitchison geometry of the simplex*<sup>5</sup>.

---

<sup>5</sup> John Aitchison. “The Statistical Analysis of Compositional Data”. In: *Journal of the Royal Statistical Society. Series B (Methodological)* 44.2 (1982), pp. 139–177, Vera Pawlowsky-Glahn, Juan José Egozcue, and Raimon Tolosana-Delgado. *Modeling and Analysis of Compositional Data*. John Wiley & Sons, 2015.

# Outline

- 1 Introduction
- 2 The Shapley values in machine learning
- 3 Compositional data analysis
- 4 Shapley composition on the simplex
- 5 Explaining a prediction with Shapley compositions
- 6 Discussion and conclusion

# The Shapley values in machine learning

We want to explain a prediction  $f(\mathbf{x})$  on the instance  $\mathbf{x} \in \mathcal{X} \subset \mathbb{R}^d$ , where  $f: \mathcal{X} \rightarrow \mathbb{R}$  is the learned model.



# The Shapley values in machine learning

We want to explain a prediction  $f(\mathbf{x})$  on the instance  $\mathbf{x} \in \mathcal{X} \subset \mathbb{R}^d$ , where  $f: \mathcal{X} \rightarrow \mathbb{R}$  is the learned model.



Some notation. Let:

- $\Pr$  be the probability distribution over  $\mathcal{X}$  of the data.

# The Shapley values in machine learning

We want to explain a prediction  $f(\mathbf{x})$  on the instance  $\mathbf{x} \in \mathcal{X} \subset \mathbb{R}^d$ , where  $f: \mathcal{X} \rightarrow \mathbb{R}$  is the learned model.



Some notation. Let:

- $\Pr$  be the probability distribution over  $\mathcal{X}$  of the data.
- $S \subseteq \mathcal{I} = \{1, 2, \dots, d\}$  be a subset of indices,

# The Shapley values in machine learning

We want to explain a prediction  $f(\mathbf{x})$  on the instance  $\mathbf{x} \in \mathcal{X} \subset \mathbb{R}^d$ , where  $f: \mathcal{X} \rightarrow \mathbb{R}$  is the learned model.



Some notation. Let:

- $\Pr$  be the probability distribution over  $\mathcal{X}$  of the data.
- $S \subseteq \mathcal{I} = \{1, 2, \dots, d\}$  be a subset of indices,
- $\mathbf{x}_S$  refers to an instance  $\mathbf{x}$  restricted to the features indicated by the indices in  $S$ .



# The Shapley values in machine learning

The **value function**:

$$\begin{aligned} v_{f,\mathbf{x},\Pr} : 2^{\mathcal{I}} &\rightarrow \mathbb{R}, \\ S &\mapsto \mathbb{E}_{\Pr}[f(\mathbf{x}) \mid \mathbf{x}_S], \end{aligned} \tag{3}$$

where  $\mathbb{E}_{\Pr}[f(\mathbf{x}) \mid \mathbf{x}_S] = \int_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x}) \Pr(\mathbf{x} \mid \mathbf{x}_S) d\mathbf{x}$ .

When an instance  $\mathbf{x}$  is observed, the expected value of the prediction is simply  $\mathbb{E}[f(\mathbf{x}) \mid \mathbf{x}] = f(\mathbf{x})$ . However, when only  $\mathbf{x}_S$  is given with  $S \neq \mathcal{I}$ , there is uncertainty about the non-observed features and we therefore compute the expected prediction given  $\mathbf{x}_S$ .

# The Shapley values in machine learning

The **contribution** of the feature indexed by  $i \notin S$  in the prediction  $f(\mathbf{x})$  given the known features indexed by  $S$  is given by:

$$c_{f,\mathbf{x},\text{Pr}}(i, \mathbf{X}_S) = v_{f,\mathbf{x},\text{Pr}}(\mathbf{X}_{S \cup \{i\}}) - v_{f,\mathbf{x},\text{Pr}}(\mathbf{X}_S), \quad (4)$$

This measures the contribution of the  $i$ th features with a particular *coalition* of features indexed by  $S$ .

# The Shapley values in machine learning

The whole contribution of the  $i$ th feature is computed by averaging this quantity over all possible coalitions of features as follows:

$$\phi_{f,\mathbf{x},\text{Pr}}(i) = \frac{1}{d!} \sum_{\pi} c_{f,\mathbf{x},\text{Pr}}(i, \pi_{\mathbf{X}}^{<i}), \quad (5)$$

where  $\pi$  is a permutation of the set  $\mathcal{I}$  of indexes and  $\pi_{\mathbf{X}}^{<i}$  is the features of  $\mathbf{X}$  coming before the  $i$ th feature in the ordering given by  $\pi$ .

# The Shapley values in machine learning

The whole contribution of the  $i$ th feature is computed by averaging this quantity over all possible coalitions of features as follows:

$$\phi_{f,\mathbf{x},\text{Pr}}(i) = \frac{1}{d!} \sum_{\pi} c_{f,\mathbf{x},\text{Pr}}(i, \pi_{\mathbf{x}}^{<i}), \quad (5)$$


where  $\pi$  is a permutation of the set  $\mathcal{I}$  of indexes and  $\pi_{\mathbf{x}}^{<i}$  is the features of  $\mathbf{X}$  coming before the  $i$ th feature in the ordering given by  $\pi$ .

This quantity is known as the **Shapley value** for the  $i$ th feature.

# The Shapley values in machine learning

It comes from cooperative game theory and is known to be the only quantity respecting a set of desired axiomatic properties<sup>6</sup>.

---

<sup>6</sup>Lloyd S Shapley et al. "A value for n-person games". In: (1953), pp 307–317. 


# The Shapley values in machine learning

It comes from cooperative game theory and is known to be the only quantity respecting a set of desired axiomatic properties<sup>6</sup>.

- Linearity with respect to the model ( $\alpha, \beta \in \mathbb{R}$ ):

$$\phi_{\alpha f + \beta g}(i) = \alpha \phi_f(i) + \beta \phi_g(i),$$

---

<sup>6</sup>Lloyd S Shapley et al. "A value for n-person games". In: (1953), pp. 307–317. 

# The Shapley values in machine learning


It comes from cooperative game theory and is known to be the only quantity respecting a set of desired axiomatic properties<sup>6</sup>.

- Linearity with respect to the model ( $\alpha, \beta \in \mathbb{R}$ ):  
 $\phi_{\alpha f + \beta g}(i) = \alpha \phi_f(i) + \beta \phi_g(i),$
- The “centered” learned model is additively separable with respect to the Shapley values:

$$f(\mathbf{x}) - \mathbb{E}_{\mathbf{P}_f}[f(\mathbf{X})] = \sum_{i=1}^d \phi_f(i), \quad (6)$$

which is known as the *efficiency* property.

---

<sup>6</sup>Lloyd S Shapley et al. “A value for n-person games”. In: (1953), pp. 307–317. 

# The Shapley values in machine learning

It comes from cooperative game theory and is known to be the only quantity respecting a set of desired axiomatic properties<sup>6</sup>.


- Linearity with respect to the model ( $\alpha, \beta \in \mathbb{R}$ ):  
 $\phi_{\alpha f + \beta g}(i) = \alpha \phi_f(i) + \beta \phi_g(i),$
- The “centered” learned model is additively separable with respect to the Shapley values:

$$f(\mathbf{x}) - \mathbb{E}_{\mathbf{P}_f}[f(\mathbf{X})] = \sum_{i=1}^d \phi_f(i), \quad (6)$$

which is known as the *efficiency* property.

- Symmetry

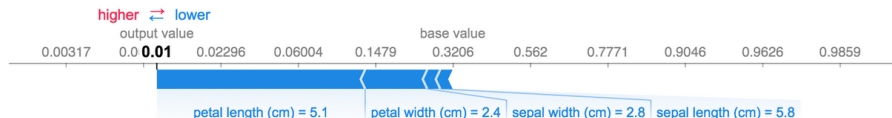
---

<sup>6</sup>Lloyd S Shapley et al. “A value for n-person games”. In: (1953), pp. 307–317. 



# The Shapley values in machine learning

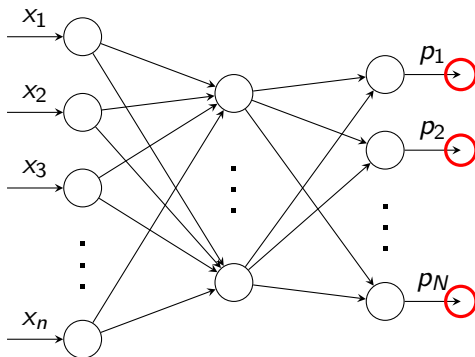
Example of explanation:



**Figure:** Explanation of the probability for the class Setosa for a flower from the Iris dataset. The classifier is an SVM with radial basis function and pairwise coupling. Image from <https://github.com/shap/shap/tree/master>.

Note that the Shapley explanation is ran in the *logit* domain!

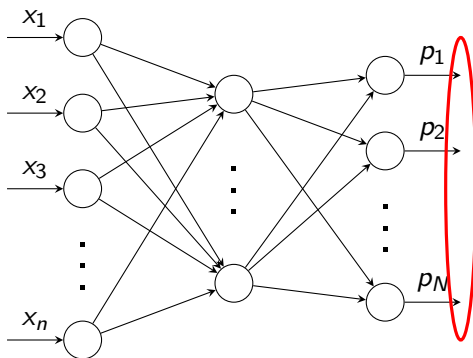
# The Shapley values in machine learning



Some explain the output one-by-one,

But a probability distribution lives on a simplex,  
**The relative information matter!!**

# The Shapley values in machine learning



We will explain the probabilities all together using the *Aitchison geometry of the simplex*.

# Compositional data analysis

# Shapley composition on the simplex

# Explaining a prediction with Shapley compositions

# Discussion and conclusion

# References I

- [1] John Aitchison. “The Statistical Analysis of Compositional Data”. In: *Journal of the Royal Statistical Society. Series B (Methodological)* 44.2 (1982), pp. 139–177.
- [2] Scott M Lundberg and Su-In Lee. “A Unified Approach to Interpreting Model Predictions”. In: *Advances in Neural Information Processing Systems 30*. Ed. by I. Guyon et al. Curran Associates, Inc., 2017, pp. 4765–4774.
- [3] Vera Pawlowsky-Glahn, Juan José Egozcue, and Raimon Tolosana-Delgado. *Modeling and Analysis of Compositional Data*. John Wiley & Sons, 2015.
- [4] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. “” Why should i trust you?” Explaining the predictions of any classifier”. In: *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*. 2016, pp. 1135–1144.



# References II

- [5] Lloyd S Shapley et al. “A value for n-person games”. In: (1953), pp. 307–317.
- [6] Erik Štrumbelj and Igor Kononenko. “Explaining prediction models and individual predictions with feature contributions”. In: *Knowledge and information systems* 41 (2014), pp. 647–665.

Thank you!!