

Explaining probabilistic predictions on the simplex with Shapley compositions

Anonymous Authors¹

Abstract

The Shapley value has been widely used for measuring the contribution of each feature to a model's prediction. However, coming from game theory, this has been designed for a one-dimensional function's codomain. For a multiclass probabilistic classifier, the output is a discrete probability distribution, over a set of more than two possible classes, and lives on a multidimensional simplex. In this case, the Shapley values are sometimes computed on each output dimension one-by-one, in an implicit one-vs-rest setting, ignoring the compositional nature of the output distribution. Indeed, elements of the simplex are known as compositional data and a discrete probability distribution can therefore be treated as such taking into account the relative information between probabilities. Using the Aitchison geometry of the simplex, this paper presents a first initiative for a multidimensional extension of the concept of Shapley value, named Shapley composition, for explaining probabilistic predictions on the simplex in machine learning.

1. Introduction

Modern machine learning approaches like the one based on deep learning are often regarded as black-boxes making them not reliable for real-life application where the machine learning prediction has to be understood. These last years, the number of contribution to make models more explainable has therefore increased in the machine learning literature. One way to better understand a prediction would be to measure the contribution of each input features on the computation of the model output. The concept of Shap-

ley value is now widely used for this purpose (Štrumbelj & Kononenko, 2014; Datta et al., 2016) especially since the release of the SHAP toolkit (Lundberg & Lee, 2017)¹. The Shapley value came from cooperative game theory...

explain shapley in game theory,

How it is applied to ML,

Limitation,

...

1.1. Contributions

...

2. The Shapley value in machine learning

In this section, we recall the theoretical formulation of the Shapley value for measuring the contribution of each feature on a machine learning prediction.

Let $f : \mathcal{X} \rightarrow \mathbb{R}$ be a learned model one want to locally explain where $f(\mathbf{x})$ is the prediction on the instance $\mathbf{x} \in \mathcal{X} \subset \mathbb{R}^d$. Let \Pr be the probability distribution over \mathcal{X} of the data². Let $S \subseteq \mathcal{I} = \{1, 2, \dots, d\}$, where d is the number of features that composes an instance $\mathbf{x} \in \mathcal{X}$, be a subset of indices. \mathbf{x}_S refers to an instance \mathbf{x} restricted to the features indicated by the indices in S .

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 other features and we therefore compute the expected prediction given \mathbf{x}_S : $\mathbb{E}_{\Pr}[f(\mathbf{x}) \mid \mathbf{x}_S] = \int_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x}) \Pr(\mathbf{x} \mid \mathbf{x}_S) d\mathbf{x}$. 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:

$$c_{f, \mathbf{x}, \Pr}(i, \mathbf{X}_S) = v_{f, \mathbf{x}, \Pr}(\mathbf{X}_{S \cup \{i\}}) - v_{f, \mathbf{x}, \Pr}(\mathbf{X}_S), \quad (1)$$

¹<https://github.com/shap/shap>

²Usually, this is unknown but expectations will be replaced by empirical averagings.

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

where v is known as the value function:

$$v_{f,\mathbf{x},\text{Pr}} : 2^{\mathcal{I}} \rightarrow \mathbb{R}, \quad (2)$$

$$S \mapsto \mathbb{E}_{\text{Pr}}[f(\mathbf{x}) \mid \mathbf{x}_S],$$

where $2^{\mathcal{I}}$ is the set of all subsets of \mathcal{I} . This measure the contribution of the i th features with a particular coalition of features indexed by S . The whole contribution of the i th feature is computing by averaging this quantity over all possible coalitions as follow:

$$\phi_{f,\mathbf{x},\text{Pr}}(i) = \frac{1}{d!} \sum_{\pi} c_{f,\mathbf{x},\text{Pr}}(i, \pi_{\mathbf{X}}^{\leq i}), \quad (3)$$

where π is a permutation of the set \mathcal{I} of indexes and $\pi_{\mathbf{X}}^{\leq i}$ is the features of \mathbf{X} coming before the i th feature in the ordering given by π . For better clarity, the subscript f,\mathbf{x},Pr will be dropt from the equations.

This quantity is known as the Shapley value for the i th feature. It comes from cooperative game theory and is known to be the only quantity respecting a set of desired axiomatic properties (Shapley et al., 1953). It is linear as a function of the model ($\alpha, \beta \in \mathbb{R}$): $\phi_{\alpha f + \beta g}(i) = \alpha \phi_f(i) + \beta \phi_g(i)$, and the “centered” learned model is additively separable with respect to the Shapley values:

$$f(\mathbf{x}) - \mathbb{E}_{\text{Pr}}[f(\mathbf{X})] = \sum_{i=1}^d \phi_f(i), \quad (4)$$

which is known as the efficiency property.

The Shapley value is designed for one-dimensinal codomain of the function f . For explaining machine learning models which output multidimensional discrete probability distribution, like in multiclass classification, people have been explaining each output dimension one-by-one, applying a logit transformation to the probabilities, resulting in a one-vs-rest comparison. However, this approach ignores the relative information between each probability and ignores the compositional nature of probability distributions. Indeed, the probabilistic output of a classifier lives on a multidimensional simplex. The latter is the sample space of data refered as compositional data we briefly review in the next section.

3. Compositional data

Compositional data carries relative information. Each element of a composition describes a part of some whole (Pawlowsky-Glahn et al., 2015) like vectors of proportions, concentrations, and discrete probability distributions. A N -part composition is a vector of N non-zero positive real numbers that sum to a constant

k . Each element of the vector is a part of the whole k . The sample space of compositional data is the simplex: $\mathcal{S}^N = \left\{ \mathbf{x} = [x_1, x_2, \dots, x_N]^T \in \mathbb{R}_+^{*N} \mid \sum_{i=1}^N x_i = k \right\}$. In a composition, only the relative information between parts matters and John Aitchison introduced the use of log-ratios of components to handle this (Aitchison, 1982). He defined several operations on the simplex which leads to what is called the Aitchison geometry of the simplex.

3.1. The Aitchison geometry of the simplex

John Aitchison defined an internal operation called perturbation, an external one called powering and an inner product (Aitchison, 2001):

- a perturbation: $\mathbf{x} \oplus \mathbf{y} = \mathcal{C}([x_1 y_1, \dots, x_N y_N])$ seen as an addition between two compositions $\mathbf{x}, \mathbf{y} \in \mathcal{S}^N$,
- a powering: $\alpha \odot \mathbf{x} = \mathcal{C}([x_1^\alpha, \dots, x_N^\alpha])$ seen as a multiplication by a scalar $\alpha \in \mathbb{R}$,
- an inner product:

$$\langle \mathbf{x}, \mathbf{y} \rangle_a = \frac{1}{2N} \sum_{i=1}^N \sum_{j=1}^N \log \frac{x_i}{x_j} \log \frac{y_i}{y_j},$$

where $\mathcal{C}(\cdot)$ is the closure operator. Since only the relative information matter, scaling factors are irrelevant and a composition \mathbf{x} is equivalent to $\lambda \mathbf{x} = [\lambda x_1, \lambda x_2, \dots, \lambda x_N]$ for all $\lambda > 0$. This equivalence is materialized by the closure operator defined for $k > 0$

as: $\mathcal{C}(\mathbf{x}) = \left[\frac{kx_1}{\|\mathbf{x}\|_1}, \frac{kx_2}{\|\mathbf{x}\|_1}, \dots, \frac{kx_N}{\|\mathbf{x}\|_1} \right]^T$, where $\mathbf{x} \in \mathbb{R}_+^{*N}$ and $\|\mathbf{x}\|_1 = \sum_{i=1}^N |x_i|$.

This give to the simplex a $(N - 1)$ -dimensional Euclidean vector space structure called Aitchison geometry of the simplex. In this paper, since we are interested in classifiers’ outputs as discrete probability distributions, we restrict ourselves to the probability simplex where $k = 1$.

3.2. The isometric log-ratio transformation

An $(N - 1)$ -dimensional orthonormal basis of the simplex, refered as an Aitchison orthonormal basis, can be built. The projection of a composition (like a discrete probability distribution) into this basis defines an isometric isomorphism between \mathcal{S}^N and \mathbb{R}^{N-1} . This is known as an Isometric-Log-Ratio (ILR) transformation (Egozcue et al., 2003) and allows to express the compositions into a Cartesian coordinates system preserving the metric of the Aitchison geometry. Within

this real space, the permutation, the powering and the Aitchison inner product defined above are respectively the standard addition, multiplication by a scalar and inner product.

Given a composition $\mathbf{p} = [p_1, \dots, p_N]^T \in \mathcal{S}^N$ we express its ILR transformation as $\tilde{\mathbf{p}} = \text{ilr}(\mathbf{p}) = [\tilde{p}_1, \dots, \tilde{p}_{N-1}]^T \in \mathbb{R}^{N-1}$. The i th element \tilde{p}_i of $\tilde{\mathbf{p}}$ is be obtained as: $\tilde{p}_i = \langle \mathbf{p}, \mathbf{e}^{(i)} \rangle_a$ where the set $\{\mathbf{e}^{(i)} \in \mathcal{S}^N, i = 1, \dots, N-1\}$ forms an Aitchison orthonormal basis of the simplex. The choice of the basis will be discussed in Section 5.2.

4. Shapley on the simplex

In Section 2 we have briefly presented the standard use of Shapley value designed for explaining one-dimensional predictions. In this Section, we will see how the Aitchison geometry can be used to extend the concept of Shapley value to the simplex for explaining multidimensional probabilistic predictions.

Let $\mathbf{f} : \mathcal{X} \rightarrow \mathcal{S}^N$ be a learned model, like a N -classes probabilistic classifier for instance, which outputs a probabilistic prediction on the $(N-1)$ -dimensional probability simplex \mathcal{S}^N . To properly consider the relative information between the probabilities, the output of the model must be treated as compositional data using the operators and metric defined by the Aitchison geometry of the simplex. We therefore rewrite the contribution and the value function of Equations 1 and 2 as follow:

$$\mathbf{c}_{\mathbf{f}, \mathbf{x}, \text{Pr}}(i, \mathbf{X}_S) = \mathbf{v}_{\mathbf{f}, \mathbf{x}, \text{Pr}}(\mathbf{X}_{S \cup \{i\}}) \ominus \mathbf{v}_{\mathbf{f}, \mathbf{x}, \text{Pr}}(\mathbf{X}_S), \quad (5)$$

where $\mathbf{a} \ominus \mathbf{b}$ is the perturbation $\mathbf{a} \oplus ((-1) \odot \mathbf{b})$ which correspond to a subtraction between compositions \mathbf{a} and \mathbf{b} , and where:

$$\mathbf{v}_{\mathbf{f}, \mathbf{x}, \text{Pr}} : 2^{\mathcal{X}} \rightarrow \mathcal{S}^N, \quad (6)$$

$$\mathbf{X}_S \mapsto \mathbb{E}_{\text{Pr}}^{\mathcal{A}}[\mathbf{f}(\mathbf{x}) \mid \mathbf{x}_S].$$

The \mathcal{A} in superscript highlight the fact that the expectation is done with respect to the Aitchison measure, rather than the Lebesgue measure, which can simply be computed as: $\mathbb{E}^{\mathcal{A}}[\mathbf{Y}] = \text{ilr}^{-1}(\mathbb{E}[\text{ilr}(\mathbf{Y})])$, where $\mathbb{E}^{\mathcal{A}}$ refers to the expectation with respect to the Aitchison measure while \mathbb{E} refers to the expectation with respect to the Lebesgue measure (Pawlowsky-Glahn et al., 2015).

The Shapley quantity expressing the contribution of the i th feature on a prediction can simply be expressed on the simplex as the Shapley composition $\phi(i)$ given by:

$$\phi_{\mathbf{f}, \mathbf{x}, \text{Pr}}(i) = \frac{1}{d!} \odot \bigoplus_{\pi} \mathbf{c}_{\mathbf{f}, \mathbf{x}, \text{Pr}}(i, \pi_{\mathbf{X}}^{\leq i}). \quad (7)$$

It can be shown (in Appendix A) that the linearity and the efficiency properties naturally hold for the Shapley composition:

$$\begin{aligned} \phi_{\alpha \odot \mathbf{f}(\mathbf{x}) \oplus \beta \odot \mathbf{g}(\mathbf{x})}(i) &= \alpha \odot \phi_{\mathbf{f}}(i) \oplus \beta \odot \phi_{\mathbf{g}}(i), \\ \bigoplus_{i=1}^d \phi_{\mathbf{f}}(i) &= \mathbf{f}(\mathbf{x}) \ominus \mathbb{E}_{\text{Pr}}^{\mathcal{A}}[\mathbf{f}(\mathbf{X})]. \end{aligned} \quad (8)$$

This can be seen as a multidimensional extension of the Shapley value framework on the simplex. Here, the Shapley quantity is not a scalar anymore, this is a composition living on the probability simplex. In the next section, we will see in more details how this can be used to explain the contribution of the features on a multidimensional probabilistic prediction.

5. Explaining probabilistic prediction with Shapley compositions

Given a prediction $\mathbf{f}(\mathbf{x})$, the Shapley composition $\phi_{\mathbf{f}, \mathbf{x}, \text{Pr}}(i)$ describes the contribution of the i th feature on the prediction. The efficiency property shows how the probability distribution moves from the base value, i.e. the expected prediction regardless of the current input, to the prediction $\mathbf{f}(\mathbf{x})$. In the standard Shapley formulation recalled in Section 2, the prediction is one-dimensional such that the Shapley quantity is a scalar. In application where there are more than two possible classes, the prediction is multidimensional such that the Shapley quantity is too. Both lives on the same space: the probability simplex. In this section, we discuss how the set of Shapley compositions can be analysed to better understand the contribution and influence of each features on the prediction.

5.1. Visualization

The Shapley compositions can be visualized in the Euclidean space isometric to the simplex thanks to the ILR transformation presented in Section 3.2. This space has the advantage of being intuitive since it is a standard real $(N-1)$ -dimensional vector space.

5.1.1. Three classes

In the three classes case, the space is 2-dimensional. We illustrate this example with the well known Iris classification dataset consisting of a set of flowers described by 4 features: sepal length and width and petal length and width. The aim is to predict to which of the three species, setosa, versicolor and virginica, a flower belongs to. In our example we use a Support Vector Machine (SVM) classifier with a radial basis function (rbf) kernel as a classifier. Figure 1 shows the explanation of one versicolor instance where 1a shows the

Shapley composition in the ILR space and 1b shows how they move the base distribution to the prediction. Having the highest norm, the petal length is the feature contributing the most on the prediction, and moves the base distribution to the versicolor maximum probability decision region (maximum probability decision region boundaries are the dashed gray lines). Being orthogonal to the virginica class-composition, this suggests that this feature does not contribute on the predicted probability for this class. The Shapley composition for the petal width goes straight to the opposite direction of the setosa class vector suggesting that this feature contributes in rejecting this class. The other Shapley compositions have a low norm suggesting these features do not contribute in the prediction.

5.1.2. Four classes

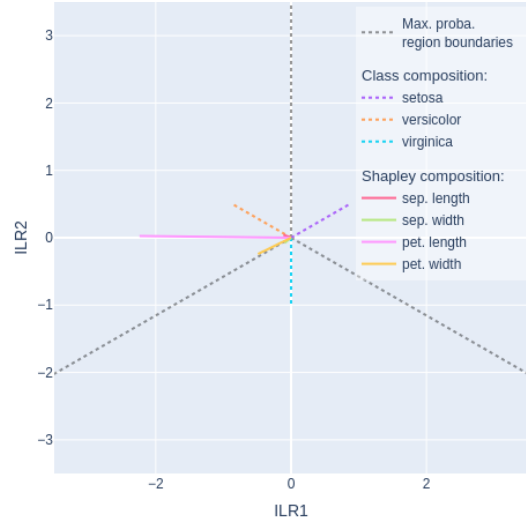
In a four classes example, the simplex is 3-dimensional. We illustrate this with a simple digit recognition task³. It consists of classifying an 8×8 image as representing one of the digit among: 0, 1, 2 and 3. Since they are 64 pixels, considering each pixel as a feature would correspond to 64 Shapley composition. We therefore reduce the number of features to 6 using a principal component analysis for better clarity and conciseness. We again use a SVM classifier with a rbf kernel. The same explanation analysis as before can be applied here but within a 3-dimensional plot as illustrated by Figure 2. To better understand how this space is divided into four regions each representing the maximum probability region for one class, one can think about the shape of a methane molecule. The hydrogens correspond to the vertices and the carbon to the center of a tetrahedron i.e. a 3-dimensional simplex. The relative position of the class-compositions in the ILR space are the same as the bonds between the carbon and a hydrogen: the angles are $\approx 109.5^\circ$. In this example, the tested instance is a 0⁴.

5.2. More classes: groups of parts and balances

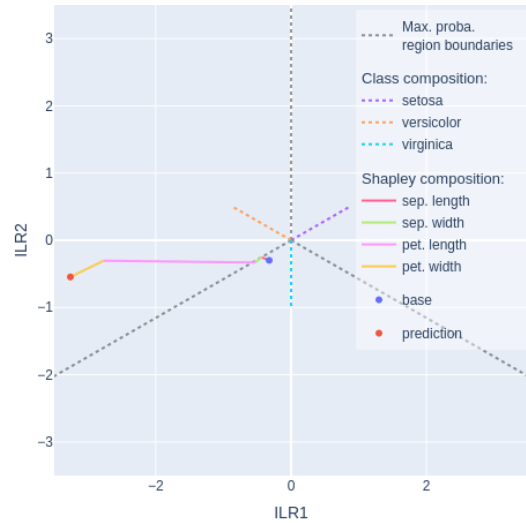
When more than three classes are involved, all the dimensions of the ILR space cannot be visualized at once, but 2 or 3-dimensional subspaces can still be visualized. In order to select the ILR components to visualize, one needs to understand what they refer to. In this section, we briefly discuss the interpretation of the ILR components.

³We use the scikit-learn's digits dataset (Pedregosa et al., 2011).

⁴More examples can be obtained from the notebooks released on the git repository: link not shared during the double blind peer reviewing.



(a) Shapley compositions in the ILR space.



(b) Sum of the Shapley compositions in the ILR space from the base distribution to the prediction.

Figure 1. Shapley explanation in the ILR space for the classification of an Iris instance.

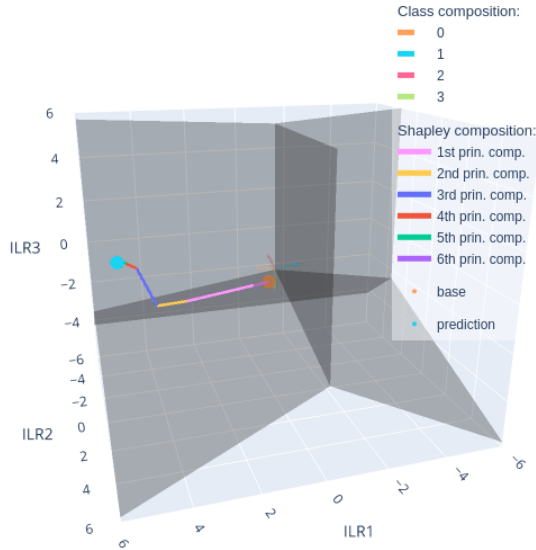


Figure 2. Shapley explanation in the 3-dimensional ILR space for a four classes digit recognition task. The Shapley compositions are summed in the ILR space from the base distribution to the prediction. The gray transparent walls markout the maximum probability decision regions.

A component of the ILR space can be interpreted as a balance, i.e. a log-ratio of two geometrical means of parts (Egozcue et al., 2003; Egozcue & Pawlowsky-Glahn, 2005; Pawlowsky-Glahn et al., 2015): one giving the central values of the probabilities in one group of classes and one for another group of classes. Therefore, a balance is here comparing the weight of two groups of classes. The set of balances is built such that they are geometrically orthogonal meaning they provide nonredundant information⁵.

This can be illustrated by a sequential binary partition or bifurcation tree. Figure 3 gives two examples: 3a shows the bifurcation tree corresponding to the basis obtained with the Gram-Schmidt procedure as in (Egozcue et al., 2003) which is the one used in the examples of Figures 1 and 2 with respectively $N = 3$ and $N = 4$. The first balance \tilde{p}_1 first compare the probabilities of class 1 and 2. Each next balance then recursively compares the probability for the next class with the probabilities for the previous ones independently of all the others.

In some applications, one may be interested in particular comparisons of groups of classes. For instance, like in an example presented in (Egozcue & Pawlowsky-Glahn, 2005), if one wants to compare political parties or groups, it may be pertinent to have a balance com-

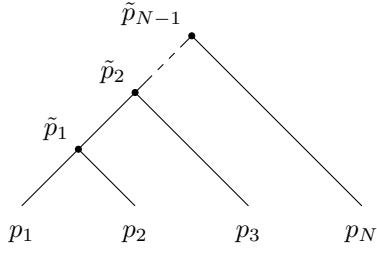
⁵Not to be confused with statistical uncorrelation (Pawlowsky-Glahn et al., 2015).

paring left and right-wing groups. However, there are sometimes no obvious relevant comparisons to study. For instance, in the handwritten digit recognition problem, it may first seem natural to compare odd with even numbers or prime with non-prime but being basically a shape recognition problem, and the shape of the numbers being independent of their arithmetic properties, these comparisons are not pertinent. Be that as it may, the choice of the basis must be left open, whether or not it is based on a relevant strategy. Let's use the basis of Figure 3b for a 10-classes digit recognition classification⁶. We comment, for conciseness, a single 2-dimensional subspace⁷. Let's have a look, in Figure 4, at the third and fifth ILR dimensions (\tilde{p}_3 and \tilde{p}_5). It is like saying we are only interested in comparing the output probability for class 0 and for class 6, and in comparing the probability for class 1 with the group of probabilities for classes 7 and 8. \tilde{p}_3 depends only on the probability for the digits 0 and 6 and \tilde{p}_5 depends only on the probabilities for the digits 1, 7 and 8. Therefore, the class-composition for the others digits have a zero projection within this subspace and are not drawn in Figure 4. The projection of the class-compositions for 0 and 6 are orthogonal to the ones for class 1, 7 and 8. Indeed, the set of classes making the balance \tilde{p}_3 and the set of classes making \tilde{p}_5 have no intersection contrary to the example of Figure 1 where \tilde{p}_1 is comparing the probabilities for the class setosa with the probability for the class versicolor and \tilde{p}_2 is comparing the probabilities for the class virginica with the group of probabilities for setosa and versicolor. In this case, none of the class-compositions are orthogonal. In Figure 4, since \tilde{p}_5 is comparing 1 with the group of digits 7 and 8, the projection of the class-compositions on this line for 1 goes in an opposite direction than the one for the class-compositions for 7 and 8. The two latters, are equal and are half as long as the former. In this way, \tilde{p}_5 compares the probability for 1 with the group of probabilities for 7 and 8 with the same weight.

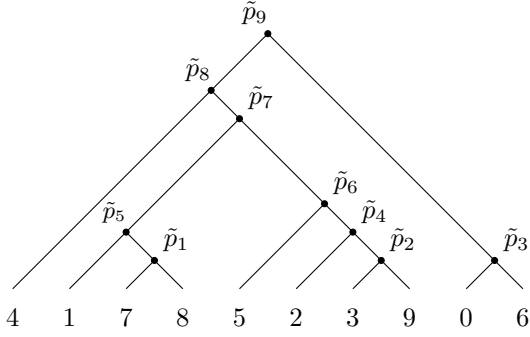
Within this space, Shapley compositions can be explored like in the examples of Figures 1 and 2 keeping in mind that this is a subset of the full ILR space.

⁶In this example, the bifurcation tree is obtained with agglomerative clustering of classes by recursively merging pair of classes based on the Mahalanobis distance in the classifier's output space, assuming that within a pair of classes, the class-conditional densities are logistic-normal (Aitchison & Shen, 1980) with same covariance matrix.

⁷More visualisations can be obtained from the shared python notebooks.



(a) Bifurcation tree corresponding to the basis obtained with the Gram-Schmidt procedure as in (Egozcue et al., 2003) and used in the examples of Figures 1 and 2.



(b) Bifurcation tree used in our 10-classes digit recognition task

Figure 3. Two examples of bifurcation tree.

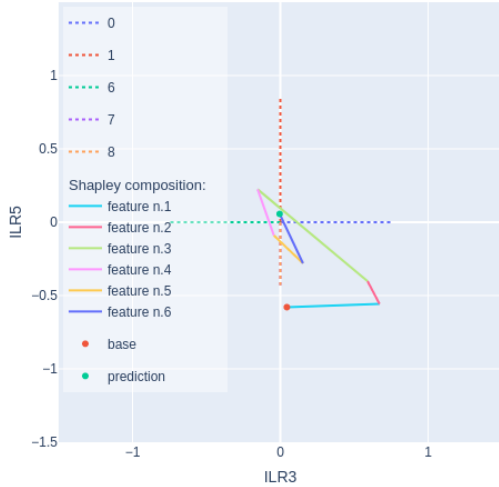


Figure 4. Sum of the Shapley compositions from the base to the prediction in the ILR subspace made of \tilde{p}_3 and \tilde{p}_5 for a test instance from the class 2. \tilde{p}_3 compares the probability for the class 0 with the one for the class 6 and \tilde{p}_5 compares the probability for the class 1 with the group of probabilities for the class 7 and 8.

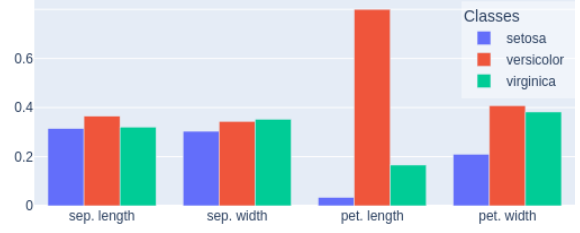


Figure 5. Shapley compositions visualized as histograms for the Iris classification example.

5.3. Angles, norms and projections

Some may find the analysis of the features contributions in cases with more than four classes tricky. Indeed, in this case, the ILR space cannot be visualized in a 2 or 3-dimensional plot and as discussed in Section 5.2, choosing which subspaces to visualize require a careful understanding of the balances and a careful building of the bifurcation tree. However, the Shapley explanation can be summarized by sets of angles, norms and projections. The norm of a Shapley composition gives the strength of the feature's contribution in the prediction. It gives the overall contribution of the feature on the prediction, regardless of its direction. The angle between two Shapley compositions can inform about their orthogonality. If two Shapley composition are orthogonal, this suggests the features are nonredundant. A negative angle would suggest that the features have an opposite influence on the prediction. The projection of a Shapley composition on the set of class-compositions informs in favor of, or against, which classes a feature is contributing. Appendix C, provides some examples of summarizing a Shapley explanation using the norm of Shapley compositions, angles between them and their projection on the class-compositions.

5.4. Histograms

If one found hard to visualize the proposed Shapley explanation in the ILR space, the Shapley composition can be visualized as histograms like discrete probability distributions. Figure 5 shows the Shapley compositions of the Iris classification example. The more uniform the histogram is, like for the sepal length and width, the less the contribution of the feature is. In opposite, the histogram for the petal length as a high value for the versicolor class, relatively to the others, confirming the contribution of the feature toward this class. For the petal width, the value for the class setosa is low relatively to the others which confirms the contribution of this feature against the class setosa.

As another illustration, Figure 6 shows the Shapley



Figure 6. Shapley compositions visualized as histograms for the seven classes digit recognition example.

compositions of the 10-classes digit recognition example. Contrary to the visualization of the compositions within the ILR space as discussed in Section 5.2, here, one can analyse all parts of each compositions within a single plot. In this example, the high of the class 2 for the first principal component confirms the contribution of this features toward this class.

5.5. About our implementation

In this work, the estimation algorithm we used to compute the Shapley compositions is an adaptation of Algorithm 2 in (Štrumbelj & Kononenko, 2014). Since the resulting Shapley compositions are approximations, the efficiency property does not necessarily hold. Each Shapley composition is therefore corrected following a similar method as in the sampling approximation in the SHAP toolkit (Lundberg & Lee, 2017)⁸ to respect the efficiency property. See Appendix D and E for more details.

6. Shapley values versus Shapley composition

We just proposed a theoretically well founded extension of the concept of Shapley value for the explaining multidimensional prediction in machine learning. The use of standard Shapley value in this context has been rarely discuss in the literature. However, some compute a Shapley value on each output dimension one-by-one. We here briefly discuss and compare this approach with the Shapley composition framework.

A Shapley value can be computed on the logit of the probability for each classes resulting in N -dimensional vector of the Shapley values for a N -classes problem. Such vector will be referred here as Shapley vector. Even if the efficiency property holds, i.e. the sum of the part-wise logit of the base distribution with the Shapley vectors for each feature is equal to the part-wise logit of the prediction, the path from the base to the prediction may go out of the simplex which is

⁸https://github.com/shap/shap/blob/master/shap/explainers/_sampling.py

counterintuitive.

7. Discussion and conclusion

Compare with standard Shapley ??

We know small expé..., sounds tricky,... first step for a theoretically founded multiclass problems explanations... ..

Features INDEPENDENCE!!

Axiomatic formulation

References

- Aitchison, J. The statistical analysis of compositional data. *Journal of the Royal Statistical Society. Series B (Methodological)*, 44(2):139–177, 1982.
- Aitchison, J. Simplicial inference. In Marlos A. G. Viana, D. S. P. R. e. A. S. S. o. A. M. i. S. (ed.), *Algebraic Methods in Statistics and Probability*, Contemporary Mathematics 287. American Mathematical Society, 2001.
- Aitchison, J. and Shen, S. M. Logistic-normal distributions: Some properties and uses. *Biometrika*, 67(2):261–272, 1980.
- Datta, A., Sen, S., and 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, 2016. doi: 10.1109/SP.2016.42.
- Egozcue, J. J. and Pawłowsky-Glahn, V. Groups of parts and their balances in compositional data analysis. *Mathematical Geology*, 37(7):795–828, 2005.
- Egozcue, J. J., Pawłowsky-Glahn, V., Mateu-Figueras, G., and Barcelo-Vidal, C. Isometric logratio transformations for compositional data analysis. *Mathematical geology*, 35(3):279–300, 2003.
- Lundberg, S. M. and Lee, S.-I. A unified approach to interpreting model predictions. In Guyon, I., Luxburg, U. V., Bengio, S., Wallach, H., Fergus, R., Vishwanathan, S., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems 30*, pp. 4765–4774. Curran Associates, Inc., 2017.
- Pawłowsky-Glahn, V., Egozcue, J. J., and Tolosana-Delgado, R. *Modeling and Analysis of Compositional Data*. John Wiley & Sons, 2015.
- Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., et al. Scikit-learn:

Machine learning in python. the Journal of machine Learning research, 12:2825–2830, 2011.

Shapley, L. S. et al. A value for n-person games. pp. 307–317, 1953.

Štrumbelj, E. and Kononenko, I. Explaining prediction models and individual predictions with feature contributions. Knowledge and information systems, 41:647–665, 2014.

A. Linearity and efficiency of the Shapley composition

In this section, we show the linearity of the Shapley composition with respect to the model prediction, and the efficiency property.

A.1. Linearity

The Shapley composition is linear, within the Aitchison geometry of the simplex, with respect to linear combination of models' predictions.

Proof. Let's consider the linear combination of predictions $\mathbf{h}(\mathbf{x}) = \alpha \odot \mathbf{f}(\mathbf{x}) \oplus \beta \odot \mathbf{g}(\mathbf{x})$. we want to check if:

$$\phi_{\mathbf{h}}(i) = \alpha \odot \phi_{\mathbf{f}}(i) \oplus \beta \odot \phi_{\mathbf{g}}(i). \quad (9)$$

We have:

$$\begin{aligned} \mathbb{E}_{\text{Pr}}^{\mathcal{A}}[\mathbf{h}(\mathbf{x}) \mid \mathbf{x}_S] &= \text{ilr}^{-1}(\mathbb{E}_{\text{Pr}}[\text{ilr}(\alpha \odot \mathbf{f}(\mathbf{x}) \oplus \beta \odot \mathbf{g}(\mathbf{x})) \mid \mathbf{x}_S]), \\ &= \text{ilr}^{-1}(\mathbb{E}_{\text{Pr}}[\alpha \text{ilr}(\mathbf{f}(\mathbf{x})) + \beta \text{ilr}(\mathbf{g}(\mathbf{x})) \mid \mathbf{x}_S]), \\ &= \text{ilr}^{-1}(\alpha \mathbb{E}_{\text{Pr}}[\text{ilr}(\mathbf{f}(\mathbf{x})) \mid \mathbf{x}_S] + \beta \mathbb{E}_{\text{Pr}}[\text{ilr}(\mathbf{g}(\mathbf{x})) \mid \mathbf{x}_S]), \\ &= \alpha \odot \text{ilr}^{-1}(\mathbb{E}_{\text{Pr}}[\text{ilr}(\mathbf{f}(\mathbf{x})) \mid \mathbf{x}_S]) \oplus \beta \odot \text{ilr}^{-1}(\mathbb{E}_{\text{Pr}}[\text{ilr}(\mathbf{g}(\mathbf{x})) \mid \mathbf{x}_S]), \\ &= \alpha \odot \mathbb{E}_{\text{Pr}}^{\mathcal{A}}[\mathbf{f}(\mathbf{x}) \mid \mathbf{x}_S] \oplus \beta \odot \mathbb{E}_{\text{Pr}}^{\mathcal{A}}[\mathbf{g}(\mathbf{x}) \mid \mathbf{x}_S]. \end{aligned} \quad (10)$$

Therefore, $\mathbf{v}_{\mathbf{h}, \mathbf{x}, \text{Pr}}(\mathbf{X}_S) = \alpha \odot \mathbf{v}_{\mathbf{f}, \mathbf{x}, \text{Pr}}(\mathbf{X}_S) \oplus \beta \odot \mathbf{v}_{\mathbf{g}, \mathbf{x}, \text{Pr}}(\mathbf{X}_S)$, meaning that \mathbf{v} is linear with respect to the learned function or model. The linearity of the contribution \mathbf{c} naturally follows:

$$\begin{aligned} \mathbf{c}_{\mathbf{h}, \mathbf{x}, \text{Pr}}(i, \mathbf{X}_S) &= \mathbf{v}_{\mathbf{h}, \mathbf{x}, \text{Pr}}(\mathbf{X}_{S \cup \{i\}}) \ominus \mathbf{v}_{\mathbf{h}, \mathbf{x}, \text{Pr}}(\mathbf{X}_S), \\ &= (\alpha \odot \mathbf{v}_{\mathbf{f}, \mathbf{x}, \text{Pr}}(\mathbf{X}_{S \cup \{i\}}) \oplus \beta \odot \mathbf{v}_{\mathbf{g}, \mathbf{x}, \text{Pr}}(\mathbf{X}_{S \cup \{i\}})) \ominus (\alpha \odot \mathbf{v}_{\mathbf{f}, \mathbf{x}, \text{Pr}}(\mathbf{X}_S) \oplus \beta \odot \mathbf{v}_{\mathbf{g}, \mathbf{x}, \text{Pr}}(\mathbf{X}_S)), \\ &= \alpha \odot \mathbf{v}_{\mathbf{f}, \mathbf{x}, \text{Pr}}(\mathbf{X}_{S \cup \{i\}}) \oplus \beta \odot \mathbf{v}_{\mathbf{g}, \mathbf{x}, \text{Pr}}(\mathbf{X}_{S \cup \{i\}}) \ominus \alpha \odot \mathbf{v}_{\mathbf{f}, \mathbf{x}, \text{Pr}}(\mathbf{X}_S) \ominus \beta \odot \mathbf{v}_{\mathbf{g}, \mathbf{x}, \text{Pr}}(\mathbf{X}_S), \\ &= \alpha \odot (\mathbf{v}_{\mathbf{f}, \mathbf{x}, \text{Pr}}(\mathbf{X}_{S \cup \{i\}}) \ominus \mathbf{v}_{\mathbf{f}, \mathbf{x}, \text{Pr}}(\mathbf{X}_S)) \oplus \beta \odot (\mathbf{v}_{\mathbf{g}, \mathbf{x}, \text{Pr}}(\mathbf{X}_{S \cup \{i\}}) \ominus \mathbf{v}_{\mathbf{g}, \mathbf{x}, \text{Pr}}(\mathbf{X}_S)), \\ &= \alpha \odot \mathbf{c}_{\mathbf{f}, \mathbf{x}, \text{Pr}}(i, \mathbf{X}_S) \oplus \beta \odot \mathbf{c}_{\mathbf{g}, \mathbf{x}, \text{Pr}}(i, \mathbf{X}_S). \end{aligned} \quad (11)$$

And the linearity of the Shap composition:

$$\begin{aligned} \phi_{\mathbf{h}}(i) &= \frac{1}{d!} \bigoplus_{\pi} \mathbf{c}_{\mathbf{h}, \mathbf{x}, \text{Pr}}(i, \pi_{\mathbf{X}}^{< i}), \\ &= \frac{1}{d!} \bigoplus_{\pi} (\alpha \odot \mathbf{c}_{\mathbf{f}, \mathbf{x}, \text{Pr}}(i, \mathbf{X}_S) \oplus \beta \odot \mathbf{c}_{\mathbf{g}, \mathbf{x}, \text{Pr}}(i, \mathbf{X}_S)), \\ &= \alpha \odot \left(\frac{1}{d!} \bigoplus_{\pi} \mathbf{c}_{\mathbf{f}, \mathbf{x}, \text{Pr}}(i, \mathbf{X}_S) \right) \oplus \beta \odot \left(\frac{1}{d!} \bigoplus_{\pi} \mathbf{c}_{\mathbf{g}, \mathbf{x}, \text{Pr}}(i, \mathbf{X}_S) \right), \\ &= \alpha \odot \phi_{\mathbf{f}}(i) \oplus \beta \odot \phi_{\mathbf{g}}(i). \end{aligned} \quad (12)$$

□

A.2. Efficiency

The efficiency property naturally holds for Shapley compositions within the Aitchison geometry.

Proof.

$$\begin{aligned}
\bigoplus_{i=1}^d \phi_{\mathbf{f}}(i) &= \bigoplus_{i=1}^d \left(\frac{1}{d!} \odot \bigoplus_{\pi} c(i, \pi_{\mathbf{X}}^{\leq i}) \right), \\
&= \frac{1}{d!} \odot \bigoplus_{i=1}^d \left(\bigoplus_{\pi} (v(\pi_{\mathbf{X}}^{\leq i+1}) \ominus v(\pi_{\mathbf{X}}^{\leq i})) \right), \\
&= \frac{1}{d!} \odot \bigoplus_{i=1}^d \left(\underbrace{\left(\bigoplus_{\pi} v(\pi_{\mathbf{X}}^{\leq i+1}) \right)}_{\mathbf{A}_{i+1}} \ominus \underbrace{\left(\bigoplus_{\pi} v(\pi_{\mathbf{X}}^{\leq i}) \right)}_{\mathbf{A}_i} \right), \\
&= \frac{1}{d!} \odot \bigoplus_{i=1}^d (\mathbf{A}_{i+1} \ominus \mathbf{A}_i), \\
&= \frac{1}{d!} \odot (\mathbf{A}_{d+1} \ominus \mathbf{A}_1), \text{ since we have a telescoping perturbation,} \\
&= \frac{1}{d!} \odot \left(\left(\bigoplus_{\pi} v(\pi_{\mathbf{X}}^{\leq d+1}) \right) \ominus \left(\bigoplus_{\pi} v(\pi_{\mathbf{X}}^{\leq 1}) \right) \right), \\
&= \frac{1}{d!} \odot \left(\left(\bigoplus_{\pi} v(\mathbf{X}) \right) \ominus \left(\bigoplus_{\pi} v(\mathbf{X}_{\emptyset}) \right) \right), \\
&= v(\mathbf{X}) \ominus v(\mathbf{X}_{\emptyset}), \text{ since } d! \text{ is the number of different permutation,} \\
&= \mathbf{f}(\mathbf{x}) \ominus \mathbb{E}_{\text{Pr}}^{\mathbf{A}}[\mathbf{f}(\mathbf{X})].
\end{aligned} \tag{13}$$

□

B. Class-compositions

A k -class-compositions $\mathbf{c}^{(k)} \in \mathcal{S}^N$ is defined as an unit norm composition going straight to the direction of the k th class. This is a discrete probability distribution with maximum probability for the k th class and uniform values for the others. The i th part of the $\mathbf{c}^{(k)}$ is:

$$c_i^{(k)} = \begin{cases} 1 - (N-1)p, & \text{if } i = k \\ p, & \text{otherwise,} \end{cases} \quad (14)$$

where $p < \frac{1}{N}$. We want the Aitchison norm of each class-composition to be one:

$$\forall k \in \{1, \dots, N\}, \quad \|\mathbf{c}^{(k)}\|_a = 1 \iff \sqrt{\frac{1}{2N} \sum_{i=1}^N \sum_{j=1}^N \left(\log \frac{c_i^{(k)}}{c_j^{(k)}} \right)^2} = 1,$$

for clarity,

we drop the (k) from the equations,

$$\begin{aligned} &\iff \sqrt{\frac{1}{2N} \sum_{i=1}^N \left((N-1) \left(\log \frac{c_i}{p} \right)^2 + \left(\log \frac{c_i}{1 - (N-1)p} \right)^2 \right)} = 1, \\ &\iff \sqrt{\frac{1}{2N} 2(N-1) \left(\log \frac{p}{1 - (N-1)p} \right)^2} = 1, \end{aligned}$$

since $p < \frac{1}{N}$ and the norm should be positive:

$$\begin{aligned} &\iff \sqrt{\frac{N-1}{N} \log \frac{1 - (N-1)p}{p}} = 1, \\ &\iff p = \frac{\exp \left(-\sqrt{\frac{N}{N-1}} \right)}{1 + (N-1) \exp \left(-\sqrt{\frac{N}{N-1}} \right)}. \end{aligned} \quad (15)$$

To summarize, the i th part of a k -class-compositions $\mathbf{c}^{(k)} \in \mathcal{S}^N$ is given by:

$$c_i^{(k)} = \frac{1}{1 + (N-1) \exp \left(-\sqrt{\frac{N}{N-1}} \right)} \left(\begin{cases} 1, & \text{if } i = k \\ \exp \left(-\sqrt{\frac{N}{N-1}} \right), & \text{otherwise,} \end{cases} \right). \quad (16)$$

In this way, $\mathbf{c}^{(k)}$ is going straight to the direction of class k and uniformly against all the others.

Table 1. Norm of the Shapley compositions, projection on the class-compositions and cosine similarity between the Shapley compositions for the Iris classification example.

	norm	projection on the class-compositions			cosine similarity between Shapley compositions			
		setosa	versicolor	virginica	pet. length	pet. width	sep. length	sep. width
pet. length	2.20	-1.91	1.90	0.02	1	.	.	.
pet. width	0.51	-0.51	0.29	0.22	0.90	1	.	.
sep. length	0.13	-0.08	0.13	-0.05	0.93	0.69	1	.
sep. width	0.11	-0.11	0.04	0.07	0.77	0.97	0.49	1

Table 2. Norm of the Shapley compositions and Cosine similarity between the Shapley composition for the 10-classes digit recognition example.

	norm	cosine similarity between Shapley compositions					
		feature n.1	feature n.2	feature n.3	feature n.4	feature n.5	feature n.6
feature n.1	4.00	1
feature n.2	1.15	-0.17	1
feature n.3	2.25	0.24	0.23	1	.	.	.
feature n.4	1.26	-0.23	0.23	0.09	1	.	.
feature n.5	0.60	0.18	-0.17	-0.01	-0.56	1	.
feature n.6	1.11	0.14	0.35	-0.19	-0.03	-0.48	1

C. Summarizing the explanation with norms, angles and projections of Shapley compositions

The Table 1 gives, for the Iris classification example of Figure 1, the norm of each Shapley compositions, the cosine similarity between them and their projection on the set of class-compositions. Having the highest norm, this confirms that the petal length is the feature that contributed the most on the prediction. Having a close to zero inner product with the virginica class-composition, shows that this feature did not contribute to the value of the probability of of this class. It contributes neither in favor, nor in the reject of this class. Having a positive inner product with the versicolor class-composition and a negative one with the setosa class-composition, this suggest that this feature is going in favor of the class versicolor and against setosa.

Table 2 shows, for the 10-classes digit recognition example, the norm of each Shapley composition and the cosine similarity between the, and Table 3 shows their projection of the set of class-compositions.

Add details...

Table 3. Projection of the Shapley compositions on the class-compositions for the 10-classes digit recognition example.

	projection on the class-compositions									
	0	1	2	3	4	5	6	7	8	9
feature n.1	1.65	-0.85	0.09	-0.57	-1.92	-0.34	1.74	1.96	0.22	-1.81
feature n.2	-0.50	0.35	0.11	0.08	-0.57	-0.29	-0.10	0.33	-0.14	0.74
feature n.3	0.36	0.45	0.23	1.20	-1.37	-0.56	0.91	-0.83	-0.31	-0.09
feature n.4	0.19	0.02	0.61	0.18	-0.53	0.46	-0.86	-0.17	-0.12	0.22
feature n.5	-0.03	-0.02	-0.22	-0.31	-0.02	0.13	0.42	-0.16	0.14	0.07
feature n.6	-0.17	0.34	-0.22	0.17	-0.12	-0.43	-0.50	0.75	0.34	-0.18

D. Estimation of the Shapley compositions

In this section, we present an adaptation of Algorithm 2 from (Štrumbelj & Kononenko, 2014) we used for the estimation the Shapley compositions.

Let d be the number of features. We want to optimally distribute the m_{\max} drawn samples over the d features. Let $\hat{\phi}_i$ be the estimation of the Shapley composition for the i th feature. We want to minimize the sum of squared errors: $\sum_{i=1}^d \|\hat{\phi}_i \ominus \phi_i\|_a^2$. Since $\hat{\phi}_i$ is a (Aitchison) sample mean we have: $\tilde{\phi}_i \approx \mathcal{N}(\tilde{\phi}_i, \frac{1}{m_i} \Sigma^{(i)})$ and $\tilde{\phi}_i - \tilde{\phi}_i \approx \mathcal{N}(\mathbf{0}, \frac{1}{m_i} \Sigma^{(i)})$ where the tilde refers to the ILR transformation. Let $\Delta_i = \tilde{\phi}_i - \tilde{\phi}_i$ and $Z_i = \|\hat{\phi}_i \ominus \phi_i\|_a = \|\tilde{\phi}_i - \tilde{\phi}_i\|_2 = \|\Delta_i\|_2$. The expectation of the sum of squared errors is:

$$\begin{aligned}
 \mathbb{E} \left[\sum_{i=1}^d Z_i^2 \right] &= \sum_{i=1}^d \mathbb{E} [Z_i^2], \\
 &= \sum_{i=1}^d \mathbb{E} \left[\sum_{j=1}^{d-1} \Delta_{ij}^2 \right], \\
 &= \sum_{i=1}^d \sum_{j=1}^{d-1} \mathbb{E} [\Delta_{ij}^2], \\
 &= \sum_{i=1}^d \sum_{j=1}^{d-1} \frac{1}{m_i} \Sigma_{jj}^{(i)}, \text{ since } \Delta_{ij} \approx \mathcal{Z}(0, \frac{1}{m_i} \Sigma_{jj}^{(i)}), \\
 &= \sum_{i=1}^d \frac{1}{m_i} \text{tr} \Sigma^{(i)}.
 \end{aligned} \tag{17}$$

When a sample is drawn, the feature for which the sample will be used for improving the Shapley composition estimation is chosen to maximize $\frac{\text{tr} \Sigma^{(i)}}{m_i} - \frac{\text{tr} \Sigma^{(i)}}{m_i+1}$. Like in (Štrumbelj & Kononenko, 2014), this is summarized in Algorithm 2.

Algorithm 1 Adaptation of the Algorithm 1 from (Štrumbelj & Kononenko, 2014) for approximating the Shapley composition of the i th feature, with model \mathbf{f} , instance $\mathbf{x} \in \mathcal{X}$ and m drawn samples.

Initialize $\phi_i \leftarrow \text{ilr}^{-1}(\mathbf{0})$

for 1 to m do

 Randomly select a permutation π of the set of indexes \mathcal{I} ,

 Randomly select a sample $\mathbf{w} \in \mathcal{X}$,

 Construct two instances:

- \mathbf{b}_1 : which takes the values from \mathbf{x} for the i th feature and the features indexed before i in the order given by π , and takes the values from \mathbf{w} otherwise,
- \mathbf{b}_2 : which takes the values from \mathbf{x} the features indexed before i in the order given by π , and takes the values from \mathbf{w} otherwise.

$\phi_i \leftarrow \phi_i \oplus \mathbf{f}(\mathbf{b}_1) \ominus \mathbf{f}(\mathbf{b}_2)$

end for

$\phi_i \leftarrow \frac{\phi_i}{m}$

Algorithm 2 Adaptation of the Algorithm 2 from (Štrumbelj & Kononenko, 2014) for approximating all the Shapley compositions by optimally distributing a maximum number of samples m_{\max} over the d features, with model \mathbf{f} , instance $\mathbf{x} \in \mathcal{X}$ and m_{\min} the minimum number of samples each feature estimation.

Initialization: $m_i \leftarrow 0, \phi_i \leftarrow \mathbf{0}, \forall i \in \{1, \dots, d\}$,
 while $\sum_{i=1}^d m_i < m_{\max}$ do
 if $\forall i, m_i \leq m_{\min}$ then
 $j = \operatorname{argmax}_i \left(\frac{\operatorname{tr} \Sigma^{(i)}}{m_i} - \frac{\operatorname{tr} \Sigma^{(i)}}{m_i+1} \right)$,
 else
 pick a j such that $m_j < m_{\min}$,
 end if
 $\phi_j \leftarrow \phi_j + \text{result of Algorithm 1 for the } j\text{th feature and } m = 1$,
 update $\Sigma^{(j)}$ using an incremental algorithm,
 $m_j \leftarrow m_j + 1$
 end while
 $\phi_i \leftarrow \frac{\phi_i}{m_i}, \forall i \in \{1, \dots, d\}$.

E. Adjustment of the estimated Shapley compositions for efficiency

In practice, the computation of the Shapley values has an exponential time complexity and we do not have necessarily access to the true distribution of the data. The Shapley values are therefore approximated using estimation algorithms like for instance the one presented in the previous appendix. However, since the obtained values are approximations, they do not necessarily respect the desired efficiency property. This point is often overlooked in the literature. In this section we write down an adjustment strategy of the estimated Shapley compositions for them to respect the efficiency property. This is a similar strategy as in the sampling approximation of the Shapley values in the SHAP toolkit (Lundberg & Lee, 2017)⁹.

Let $\{\hat{\phi}_i\}_{1 \leq i \leq d}$ be the estimated Shapley compositions (given by the Algorithm 2 in our experiments). Let $\mathbf{s}_{err} = \mathbf{f}(\mathbf{x}) \ominus \mathbf{f}_0 \ominus \bigoplus_{i=1}^d \hat{\phi}_i$, where \mathbf{f}_0 is the base distribution, be the error composition on the perturbation of all Shapley compositions, i.e. the error making the efficiency property unfulfilled. In order to respect the efficiency property, we want this error to be the neutral element of the perturbation, i.e. the “zero” in the sense of the Aitchison geometry: the uniform distribution. We could simply perturb each estimated Shapley compositions by $\frac{1}{d} \odot \mathbf{s}_{err}$ however this would move each Shapley composition by the same amount while we want to allow the compositions with a higher estimation variance (i.e. with a precision likely to be lower) to move more than the ones with a smaller variance (i.e. with a precision likely to be higher).

We therefore weight the i th adjustment by a scalar $w_i = w(\operatorname{tr}(\Sigma^{(i)}))$, where w is an increasing function, and where $\sum_{i=1}^d w_i = 1$. Note that the vector of weight is actually a composition too. Similarly to the SHAP toolkit implementation, we choose w as:

$$w_i = w\left(\operatorname{tr}(\Sigma^{(i)})\right) = \frac{v_i}{1 + \sum_{j=1}^d v_j}, \text{ where } v_i = \frac{\operatorname{tr}(\Sigma^{(i)})}{\epsilon \max_j \operatorname{tr}(\Sigma^{(j)})}. \quad (18)$$

The i th estimated Shapley composition is then adjusted as follow:

$$\hat{\phi}_i \leftarrow \hat{\phi}_i \oplus (w_i \odot \mathbf{s}_{err}). \quad (19)$$

⁹https://github.com/shap/shap/blob/master/shap/explainers/_sampling.py

In this way, when ϵ goes to zero¹⁰, the efficiency property is respected for the adjusted Shapley compositions and more weight is given to the adjustments of the Shapley compositions with a higher estimation variance.

¹⁰In our experiments, $\epsilon = 10^{-6}$.