
Supplementary Materials :

Accurate and robust Shapley Values for explaining predictions and identifying group of Important Variables

Contents

A	Proofs	2
B	Fast Algorithm for the computation of Shapley Values with the Leaf estimator	3
C	Same Decision Probability with tree-based model	6
D	Link between the Algorithm 1 (TreeSHAP with path-dependent) and \hat{f}^{SHAP}	7
E	Focus on influential variables in a linear model	9
F	Additional examples	10
F.1	Impact of quantile discretization	10
F.2	The differences between Coalition and sum on Census Data	10
F.3	SDP and Active SV analysis on Lucas Data	11
G	Individual Shapley values for indicator variables in Dummy Encoding	12
H	Plug-In estimator of Marginal expectation	14
I	EXPERIMENTAL SETTINGS	15
A.1	Toy model of Section 2.3	15
A.2	Toy model of Section 3.3	15
A.3	Comparisons of calculation time of SV estimates with ACV and TreeSHAP	15

A Proofs

This section gathers all the proofs of the propositions and claims of the main paper.

2. Coalition and Invariance for Shapley Values

2.1 Invariance under reparametrization for continuous variables

Proposition A.1. *Let f and $\tilde{f} = f \circ \varphi^{(-1)}$ its reparametrization, then we have for all $i \in \llbracket 1, p \rrbracket$, for all $\mathbf{x}, \mathbf{u} = \varphi(\mathbf{x})$:*

$$\phi_i(f, \mathbf{x}) = \phi_i(\tilde{f}, \varphi(\mathbf{x})).$$

Proof. It is a direct application of the change of variables formula. If $g(\mathbf{x})$ is the joint density of X_1, \dots, X_p (X_i has density g_i), the transformed variable $\mathbf{U} = (\varphi_1(X_1), \dots, \varphi_p(X_p))$ has density $\tilde{g}(\mathbf{u}) = g(\varphi^{(-1)}\mathbf{u}) \times \prod_i |J(\varphi_i^{(-1)})(u_i)|$. With obvious notations, we have

$$\tilde{g}(u_{\bar{S}}|u_S) = \frac{\tilde{g}(u_{\bar{S}}, u_S)}{\tilde{g}_S(u_S)} = g\left(\varphi_{\bar{S}}^{(-1)}(u_{\bar{S}}|\varphi_S^{(-1)}(u_S))\right) \times \prod_{i \in \bar{S}} |J(\varphi_i^{(-1)})(u_i)|.$$

The computation of the reduced predictor is straightforward

$$\begin{aligned} E[f(\mathbf{X})|\mathbf{x}_S] &= \int f(\mathbf{x}_S, \mathbf{x}_{\bar{S}})g(\mathbf{x}_{\bar{S}}|\mathbf{x}_S)d\mathbf{x}_{\bar{S}} \\ &= \int f(\varphi_S^{(-1)}(\varphi_S(\mathbf{x}_S)), \varphi_{\bar{S}}^{(-1)}(\varphi_{\bar{S}}(\mathbf{x}_{\bar{S}})))g(\mathbf{x}_{\bar{S}}|\mathbf{x}_S)d\mathbf{x}_{\bar{S}} \\ &= \int \tilde{f}(\mathbf{u}_S, \mathbf{u}_{\bar{S}})g\left(\varphi_{\bar{S}}^{(-1)}(\mathbf{u}_{\bar{S}}|\varphi_S^{(-1)}(\mathbf{u}_S))\right) \prod_{i \in \bar{S}} |J(\varphi_i^{(-1)})(u_i)| d\mathbf{u}_{\bar{S}} \\ &= E\left[\tilde{f}(\mathbf{U}_S, \mathbf{U}_{\bar{S}})|\mathbf{U}_S = \mathbf{u}_S\right]. \end{aligned}$$

The equality of Shapley Values is then a direct consequence of the equality of reduced predictors. \square

2.2 Invariance for encoded categorical variable

We recall the expression of the SV for 2 variables for all $x \in \mathbb{R}$ and $Y \in \{1, \dots, K\}$. The role of variable X, Y are symmetric and the categorical or quantitative nature of the variable does not have any impact on the computation of SV given:

$$\begin{cases} \phi_X(f; x, y) = \frac{1}{2} (E[f(X, Y)|X = x] - E[f(X, Y)]) + \frac{1}{2} (f(x, y) - E[f(X, Y)|Y = y]) \\ \phi_Y(f; x, y) = \frac{1}{2} (E[f(X, Y)|Y = y] - E[f(X, Y)]) + \frac{1}{2} (f(x, y) - E[f(X, Y)|X = x]) \end{cases} \quad (\text{A.1})$$

Proposition A.2. *For all $x \in \mathcal{X}$, and if $y_{1:K-1} = \mathcal{C}(y)$ then*

$$\begin{cases} \phi_C(\tilde{f}; x, y_{1:K-1}) &= \phi_Y(f; x, y) \\ \phi_X(\tilde{f}; x, y_{1:K-1}) &= \phi_X(f; x, y) \end{cases} \quad (\text{A.2})$$

Proof. As we consider only doable $(x, y_{1:K-1})$, then $\exists! y \in \{1, \dots, K\}$ such that $\mathcal{C}(y) = y_{1:K-1}$. We have the coalition $C = \{1, \dots, K-1\}$, and number of variables $p = K$, meaning

$$\phi_{\{1, \dots, K-1\}}(\tilde{f}; x, y_{1:K-1}) = \frac{1}{2} \left\{ \frac{1}{\binom{1}{0}} \Delta(\tilde{f}; \emptyset, C) + \frac{1}{\binom{1}{1}} \Delta(\tilde{f}; \{X\}, C) \right\}$$

where

$$\begin{aligned} \Delta(\tilde{f}; \emptyset, C) &= E_{\tilde{P}} \left[\tilde{f}(X, Y_{1:K-1}) | Y_{1:K-1} = y_{1:K-1} \right] - E_{\tilde{P}} \left[\tilde{f}(X, Y_{1:K-1}) | \emptyset \right] \\ &= E_P \left[\tilde{f}(X, \varphi(Y)) | Y = y \right] - E_P \left[\tilde{f}(X, \varphi(Y)) \right] \\ &= E_P \left[f(X, Y) | Y = y \right] - E_P \left[f(X, Y) \right] \end{aligned}$$

Indeed

$$\begin{aligned}
E_{\tilde{P}} \left[\tilde{f}(X, Y_{1:K-1}) | Y_{1:K-1} = y_{1:K-1} \right] &= \int \tilde{f}(x, y_{1:K-1}) dP(x | y_{1:K-1}) \\
&= \int \tilde{f}(x, y_{1:K-1}) \frac{dP(x, y_{1:K-1})}{P(y_{1:K-1})} \\
&= \int \tilde{f}(x, \varphi(y)) \frac{dP(x, \varphi(y))}{P(\varphi(y))} \\
&= \int f(x, y) \frac{dP(x, y)}{P(y)}
\end{aligned}$$

In addition,

$$\begin{aligned}
\Delta(\tilde{f}; \{X\}, C) &= E_{\tilde{P}} \left[\tilde{f}(X, Y_{1:K-1}) | X = x, Y_{1:K-1} = y_{1:K-1} \right] - E_{\tilde{P}} \left[\tilde{f}(X, Y_{1:K-1}) | X = x \right] \\
&= \tilde{f}(x, y_{1:K-1}) - E_P \left[\tilde{f}(X, \varphi(Y)) | X = x \right] \\
&= \tilde{f}(x, \varphi(y)) - E_P \left[\tilde{f}(X, \varphi(Y)) | X = x \right] \\
&= f(x, y) - E_P [f(X, y) | X = x]
\end{aligned}$$

$$\begin{aligned}
\phi_{\{1, \dots, K-1\}}(\tilde{f}; x, y_{1:K-1}) &= \frac{1}{2} (E_P [f(X, Y) | Y = y] - E_P [f(X, Y)]) \\
&\quad + \frac{1}{2} (f(x, y) - E_P [f(X, y) | X = x])
\end{aligned}$$

We can recognize that we have exactly $\phi_{\{1, \dots, K-1\}}(\tilde{f}; x, y_{1:K-1}) = \phi_Y(f; x, y)$. From Equation 2.1, we derive that $\phi_X(\tilde{f}; x, y_{1:K-1}) = \phi_X(f; x, y)$. \square

Proposition A.3. *If $X \sim \mathcal{N}(\mu, \Sigma)$, then $X_{\bar{S}} | X_S = x_S$ is also multivariate gaussian with mean $\mu_{\bar{S}|S}$ and covariance matrix $\Sigma_{\bar{S}|S}$ equal:*

$$\mu_{\bar{S}|S} = \mu_{\bar{S}} + \Sigma_{\bar{S}, S} \Sigma_{S, S}^{-1} (x_S - \mu_S) \text{ and } \Sigma_{\bar{S}|S} = \Sigma_{\bar{S} \bar{S}} - \Sigma_{\bar{S} S} \Sigma_{S S}^{-1} \Sigma_{S, \bar{S}}$$

B Fast Algorithm for the computation of Shapley Values with the Leaf estimator

In section 3.2. of the main paper, we have introduced a plug-in estimator of the conditional expectation

$$f_S(\mathbf{x}_S) = \sum_{m=1}^M f_m P_X(L_m | \mathbf{X}_S = \mathbf{x}_S)$$

that is based on an approximation of the conditional expectation by event $\{\mathbf{X}_S = \mathbf{x}_S\}$ by a conditional expectation based on event $\{\mathbf{X}_S \in L_m^S\}$. For sake of notational simplicity, we write simply $L_m^S = L_m^S(\mathbf{x})$ and we remove the dependence on \mathbf{x} .

Thanks to this approximation, we can propose a straightforward estimate based on empirical frequencies, and we focus here on the computational efficiency offered by this approximation. It is well-known that the complexity of the computation of a Shapley value is exponential as we need to compute 2^p different coalitions for each observation \mathbf{x} . We show below that the complexity can be made much lower, as we derive an algorithm with complexity exponential in the depth of the tree instead of being exponential in the total number of variable p . This is very interesting as the depth of the tree is rarely above 10 in practice, while p can be very large (different order of magnitudes). We want to compute the predictor

$$\tilde{f}_S(\mathbf{x}_S) = \sum_{m=1}^M f_m P_X(L_m | \mathbf{X}_S \in L_m^S(\mathbf{x}))$$

(or its estimated version equal to $\hat{f}_S^{(Leaf)}(\mathbf{x}_S)$) that can be used for defining a new cooperative game based on the value function

$$S \mapsto \tilde{v}(f, S) \triangleq \tilde{f}_S(\mathbf{x}_S).$$

For any coalition C , our estimate of the Shapley value $\phi_C(\mathbf{x})$ is the Shapley value of the cooperative game $\tilde{v}(f, S)$ defined as

$$\tilde{\phi}_C(\tilde{f}; \mathbf{x}) = \frac{1}{p - |C| + 1} \sum_{k=0}^{p-|C|} \frac{1}{\binom{p-|C|}{k}} \sum_{S \in \mathcal{S}_k(C)} \left(\tilde{f}_{S \cup C}(\mathbf{x}_{S \cup C}) - \tilde{f}_S(\mathbf{x}_S) \right) \quad (\text{B.1})$$

We show in the next proposition that the game \tilde{v} can be split into the sum of smaller games (we consider only $C = \{i\}$ in the proposition, but it remains true for any coalition C).

Proposition B.1. *Let $f(\mathbf{x}) = \sum_{m=1}^M f_m \mathbb{1}_{L_m}(\mathbf{x})$ be a tree based on p variables $\mathbf{X} = (X_1, \dots, X_p)$. We introduce for each leaf L_m the set of variables $S_m = \{X_{N_1}, X_{N_2}, \dots, X_{N_{d_m}}\}$ used in the tree path defining the leaf L_m . For any variable X_i , the SV $\phi_i(\tilde{f}, \mathbf{x})$ can be decomposed into the sum of M cooperative games defined on each leaf L_m , and we have*

$$\tilde{\phi}_i(\tilde{f}, \mathbf{x}) = \sum_{m=1}^M \tilde{\phi}_i^m(\tilde{f}, \mathbf{x})$$

where $\tilde{\phi}_i^m(\tilde{f}, \mathbf{x})$ is a reweighted version of the Shapley Value of the cooperative game with value function

$$\tilde{v}(\tilde{f}, S) = P_X(L_m | \mathbf{X}_S \in L_m^S(\mathbf{x}))$$

Proof. By definition, we have for a single variable

$$\begin{aligned} \tilde{\phi}_i(\mathbf{x}) &= \frac{1}{p} \sum_{S \subseteq [p] \setminus \{i\}} \binom{p-1}{|S|}^{-1} \left(\tilde{f}_{S \cup i}(\mathbf{x}_{S \cup i}) - \tilde{f}_S(\mathbf{x}_S) \right) \\ &= \frac{1}{p} \sum_{S \subseteq [p] \setminus \{i\}} \binom{p-1}{|S|}^{-1} \left(\sum_{m=1}^M f_m \left[P(L_m | \mathbf{X}_{S \cup i} \in L_m^{S \cup i}) - P(L_m | \mathbf{X}_S \in L_m^S) \right] \right) \\ &= \frac{1}{p} \sum_{m=1}^M \sum_{S' \subseteq S_m \setminus \{i\}} \left[\binom{p-1}{|S'|}^{-1} f_m \left[P(L_m | \mathbf{X}_{S' \cup i} \in L_m^{S' \cup i}) - P(L_m | \mathbf{X}_{S'} \in L_m^{S'}) \right] \right. \\ &\quad \left. + \sum_{Z \neq \emptyset, Z \subseteq \overline{S_m \cup i}} \binom{p-1}{|Z| + |S'|}^{-1} f_m \left[P(L_m | \mathbf{X}_{S' \cup Z \cup i} \in L_m^{S' \cup Z \cup i}) - P(L_m | \mathbf{X}_{S' \cup Z} \in L_m^{S_m \cup Z}) \right] \right] \end{aligned}$$

However, if $Z \subseteq \bar{S}_m, S \subseteq S_m$:

$$P_X(L_m | \mathbf{X}_{Z \cup S} \in L_m^{Z \cup S}) = P_X(L_m | \mathbf{X}_S \in L_m^S). \quad (\text{B.2})$$

We shall remark that the identity of eq. (B.2) is not true anymore if we consider the conditional probability $\mathbf{X}_S = \mathbf{x}_x$. Therefore, the SV $\tilde{\phi}_i(\mathbf{x})$ can be rewrite as:

$$\begin{aligned} \tilde{\phi}_i(\mathbf{x}) &= \frac{1}{p} \sum_{m=1}^M \sum_{S' \subseteq S_m \setminus \{i\}} \left[\binom{p-1}{|S'|}^{-1} f_m \left[P(L_m | \mathbf{X}_{S' \cup i} \in L_m^{S' \cup i}) - P(L_m | \mathbf{X}_{S'} \in L_m^{S'}) \right] \right. \\ &\quad \left. + \sum_{Z \neq \emptyset, Z \subseteq \overline{S_m \cup i}} \binom{p-1}{|Z| + |S'|}^{-1} f_m \left[P(L_m | \mathbf{X}_{S' \cup i} \in L_m^{S' \cup i}) - P(L_m | \mathbf{X}_{S'} \in L_m^{S'}) \right] \right] \\ &= \frac{1}{p} \sum_{m=1}^M \sum_{S' \subseteq S_m \setminus \{i\}} \left[\binom{p-1}{|S'|}^{-1} + \sum_{Z \neq \emptyset, Z \subseteq \overline{S_m \cup i}} \binom{p-1}{|Z| + |S'|}^{-1} \right] f_m \left[P(L_m | \mathbf{X}_{S' \cup i} \in L_m^{S' \cup i}) - P(L_m | \mathbf{X}_{S'} \in L_m^{S'}) \right] \\ &\triangleq \sum_{m=1}^M \tilde{\phi}_i^m(\mathbf{x}) \end{aligned}$$

$$\sum_{Z \neq \emptyset, Z \subseteq \overline{S_m \cup i}} \binom{p-1}{|Z| + |S'|}^{-1}.$$

The algorithm is describes below, we use the following notation $\mathbb{1}_{L_m^\emptyset}(\mathbf{x}_\emptyset) = 1$.

Algorithm 1: *Multi-Games Algorithm* - Compute SV in the worst case in $\mathcal{O}(p \times M \times 2^{\text{tree-depth}})$

Inputs: $x, f(x) = \sum_{m=0}^M f_m \mathbb{1}_{L_m}(x);$

$$p = \text{length}(\mathbf{x});$$
$$\phi = \text{zeros}(p);$$
for $m = 1$ *to* M **do****for** i *in* $[p]$ **do****if i not in S_m then**

```
continue ;
```

```
/* skip to next variable */
```

end

for $S \subseteq S_m$ **do**
$$\phi[i]_+ = \left(\binom{p-1}{|S|}^{-1} + \sum_{k=1}^{p-|S_m|} \binom{p-|S_m|}{k} \binom{p-1}{k+|S|}^{-1} \right) \times$$
$$\left(\mathbb{1}_{L_m^{S_{\cup i}}}(\mathbf{x}_{S_{\cup i}}) \hat{P}_X^{(Leaf)}(L_m | \mathbf{X}_S \in L_m^{S_{\cup i}}) - \mathbb{1}_{L_m^S}(\mathbf{x}_S) \hat{P}_X^{(Leaf)}(L_m | \mathbf{X}_S \in L_m^S) \right)$$

end

end

end

return ϕ

Remark B.1. The algorithm can be vectorized in order to compute SV of several observations at the same time.

C Same Decision Probability with tree-based model

Same Decision Probability

Our methodology for identifying the most important features is based on the Same Decision Probability (SDP) criterion, introduced in [2]. In particular, we derive an efficient way to approximate the SDP for tree-based classifier.

Definition C.1. (Same Decision Probability of a classifier). Let $f : \mathcal{X} \rightarrow [0, 1]$ a probabilistic predictor and its classifier $C(\mathbf{x}) = \mathbb{1}_{f(\mathbf{x}) \geq T}$ with threshold T , the Same Decision Probability of coalition $S \subset \llbracket 1, p \rrbracket$, w.r.t $\mathbf{x} = (\mathbf{x}_S, \mathbf{x}_{\bar{S}})$ is

$$SDP_S(C; \mathbf{x}) = P(C(\mathbf{x}_S, \mathbf{X}_{\bar{S}}) = C(\mathbf{x}) | \mathbf{X}_S = \mathbf{x}_S)$$

SDP gives the probability to keep the same decision $C(\mathbf{x})$ when we do not observe the variables $\mathbf{X}_{\bar{S}}$. The higher is the probability, the better is the explanation based on S . Therefore, we want to identify the **minimal** subset of features such that the classifier makes the same decision with high probability π , given only them. More formally:

Definition C.2. (Sufficient Coalition). Given C a binary classifier, an observation $\mathbf{x} = (\mathbf{x}_S, \mathbf{x}_{\bar{S}})$, $S \triangleq S_\pi^*(\mathbf{x})$ is a Sufficient Coalition for probability π if:

1. $SDP_{S_\pi^*(\mathbf{x})}(C; \mathbf{x}) \geq \pi$
2. No subset Z of $S_\pi^*(\mathbf{x})$ satisfies $SDP_Z(f; \mathbf{x}) \geq \pi$

In order to find the coalition $S_\pi^*(\mathbf{x})$, we need to be able to compute the SDP for any subset S . However, computing the SDP is known to be computationally hard. Even for a simple Naive Bayes model and classifier, computing SDP is NP-hard [3]. Consequently, approximate criterion based on expectations instead of probabilities have been introduced see [7]. In that section, we show that we can compute exactly and efficiently the Same Decision Probability in tree-based model by relying on reduced predictors.

Proposition C.1. Let f a probabilistic predictor and its binary classifier C with threshold T , $Q_{S, \mathbf{x}}$ the law of $\mathbf{X}_{\bar{S}} | \mathbf{X}_S = \mathbf{x}_S$, then the $SDP_S(f; \mathbf{x})$ can be written explicitly with the reduced predictor:

$$SDP_S(f; \mathbf{x}) = \frac{\mathbb{E}_{\mathbf{X}_{\bar{S}} | \mathbf{X}_S = \mathbf{x}_S} [f(\mathbf{x}_S, \mathbf{X}_{\bar{S}})] - \mathbb{E}_{\mathbf{X}_{\bar{S}} | \mathbf{X}_S = \mathbf{x}_S} [f(\mathbf{x}_S, \mathbf{X}_{\bar{S}}) | f(\mathbf{x}_S, \mathbf{X}_{\bar{S}}) < T]}{\mathbb{E}_{\mathbf{X}_{\bar{S}} | \mathbf{X}_S = \mathbf{x}_S} [f(\mathbf{x}_S, \mathbf{X}_{\bar{S}}) | f(\mathbf{x}_S, \mathbf{X}_{\bar{S}}) \geq T] - \mathbb{E}_{\mathbf{X}_{\bar{S}} | \mathbf{X}_S = \mathbf{x}_S} [f(\mathbf{x}_S, \mathbf{X}_{\bar{S}}) | f(\mathbf{x}_S, \mathbf{X}_{\bar{S}}) < T]} \quad (\text{C.1})$$

Proof. First note that $SDP_S(f; \mathbf{x}) = \mathbb{P}_{\mathbf{X}_{\bar{S}} | \mathbf{X}_S = \mathbf{x}_S} [f(\mathbf{x}_S, \mathbf{X}_{\bar{S}}) \geq T]$

$$\begin{aligned} \mathbb{E}_{\mathbf{X}_{\bar{S}} | \mathbf{X}_S = \mathbf{x}_S} [f(\mathbf{x}_S, \mathbf{X}_{\bar{S}})] &= \mathbb{E}_{\mathbf{X}_{\bar{S}} | \mathbf{X}_S = \mathbf{x}_S} [f(\mathbf{x}_S, \mathbf{X}_{\bar{S}}) | f(\mathbf{x}_S, \mathbf{X}_{\bar{S}}) < T] \mathbb{P}_{\mathbf{X}_{\bar{S}} | \mathbf{X}_S = \mathbf{x}_S} [f(\mathbf{x}_S, \mathbf{X}_{\bar{S}}) < T] \\ &\quad + \mathbb{E}_{\mathbf{X}_{\bar{S}} | \mathbf{X}_S = \mathbf{x}_S} [f(\mathbf{x}_S, \mathbf{X}_{\bar{S}}) | f(\mathbf{x}_S, \mathbf{X}_{\bar{S}}) \geq T] \mathbb{P}_{\mathbf{X}_{\bar{S}} | \mathbf{X}_S = \mathbf{x}_S} [f(\mathbf{x}_S, \mathbf{X}_{\bar{S}}) \geq T] \end{aligned}$$

Rearranging the terms leads to equation C.1 □

Based on the computation of the SDP of any coalition given by the previous propositions, we can derive an algorithm that finds the Sufficient Coalitions for probability π i.e $S_\pi^*(\mathbf{x})$.

Unlike SV computation, we don't have to compute all the conditional expectations for all subsets in order to find the coalition S_π^* . We use a greedy algorithm that computes the SDPs for subsets of increasing sizes (starting from 1) until we find a minimal subset satisfying the Sufficient Coalition conditions. The algorithm is described in 2 and defines the function *returnSubsets*(\mathbf{x} , *size*) that returns all subsets of length *size* of \mathbf{x} .

We already know how to estimate $E[f(\mathbf{x}_S, \mathbf{X}_{\bar{S}}) | \mathbf{X}_S = \mathbf{x}_S]$. Therefore, we use the same idea to estimate $\mathbb{E}_{\mathbf{X}_{\bar{S}} | \mathbf{X}_S = \mathbf{x}_S} [f(\mathbf{x}_S, \mathbf{X}_{\bar{S}}) | f(\mathbf{x}_S, \mathbf{X}_{\bar{S}}) < T]$. We estimate each probability of the Leaf estimator of the reduced predictor as

$$\hat{P}_{\mathbf{X}_{\bar{S}} | \mathbf{X}_S = \mathbf{x}_S}^{(Leaf)} [f(\mathbf{x}_S, \mathbf{X}_{\bar{S}}) | f(\mathbf{x}_S, \mathbf{X}_{\bar{S}}) < T] = \frac{N(L_m, f < T)}{N(L_m^S, f < T)}$$

and for the Discrete estimator, we use

$$\hat{P}_{\mathbf{X}_{\bar{S}}|\mathbf{X}_S=\mathbf{x}_S}^{(D)} [f(\mathbf{x}_S, \mathbf{X}_{\bar{S}}) | f(\mathbf{x}_S, \mathbf{X}_{\bar{S}}) < T] = \frac{N(L_m, \mathbf{x}_S, f < T)}{N(\mathbf{x}_S, f < T)}$$

where

- $N(\mathbf{x}_S, f < T)$: number of observations such that $\mathbf{X}_S = \mathbf{x}_S$ and $f(x) < T$ (across all the leaves of the tree)
- $N(L_m, \mathbf{x}_S, f < T)$: number of observations in leaf L_m that meet the condition $\mathbf{X}_S = \mathbf{x}_S$ and $f(x) < T$
- $N(L_m, f < T)$: number of observations in the leaf L_m that meet the condition $f(x) < T$,
- $N(L_m^S, f < T)$: number of observations that meet the conditions $\mathbf{x}_S \in L_m^S$ and $f(x) < T$ across all the leaves of the tree.

Algorithm 2: Find Sufficient Coalition

Inputs: \mathbf{x}, π ;
 $n = \text{length}(\mathbf{x})$;
 $\text{find} = \text{False}$;
 $\text{bestSdp} = -1$;
for $\text{size} = 1$ **to** n **do**
 for $S \subset \text{returnSubsets}(\mathbf{x}, \text{size})$ **do**
 $\text{sdp} = \text{SDP}_S(\mathbf{x}, f)$;
 if $\text{sdp} \geq \pi$ **and** $\geq \text{bestSdp}$ **then**
 $\text{bestSdp} = \text{sdp}$;
 $S_\pi^* = S$;
 $\text{find} = \text{True}$;
 end
 end
end
if $\text{find} == \text{True}$ **then**
 return S_π^*
end

D Link between the Algorithm 1 (TreeSHAP with path-dependent) and \hat{f}^{SHAP}

In section 3.1, we have said that the recursive algorithm 1 introduced in [6] and shows in figure 2 assumes that the probabilities can be factored with the decision tree as:

$$P_X^{SHAP} \left(\prod_{k=1}^{d_m} I_{N_k} | \mathbf{X}_S = \mathbf{x}_S \right) = \prod_{i=2|N_i \notin S}^{d_m} P(X_{N_i} \in I_{N_i} | X_{N_{i-1}} \in I_{N_{i-1}}) \times \delta_S(N_1) \quad (\text{D.1})$$

with $\delta_S(N_1) = P(X_{N_1} \in I_{N_1})$ if $N_1 \notin S$, and 1 otherwise.

To show the link between \hat{f}^{SHAP} and the Algorithm 1, let choose an observation $x_{ref} = [2, 3, 0.5, -1]$ and compute $E[f(\mathbf{X}) | X_0 = 2, X_2 = 0.5]$ where f is the tree in figure 1. The condition $X_0 = 2, X_2 = 0.5$ is compatible with the leaves 6, 7, 11, 13, 14, we denote $f_6, f_7, f_{11}, f_{13}, f_{14}$ the value of each leaf respectively. The output of the algorithm 1 (described in figure 2) is on step 4 of the table D.1 below and its corresponds to $\hat{f}_{\{0,2\}}^{(SHAP)}(x_{ref})$.

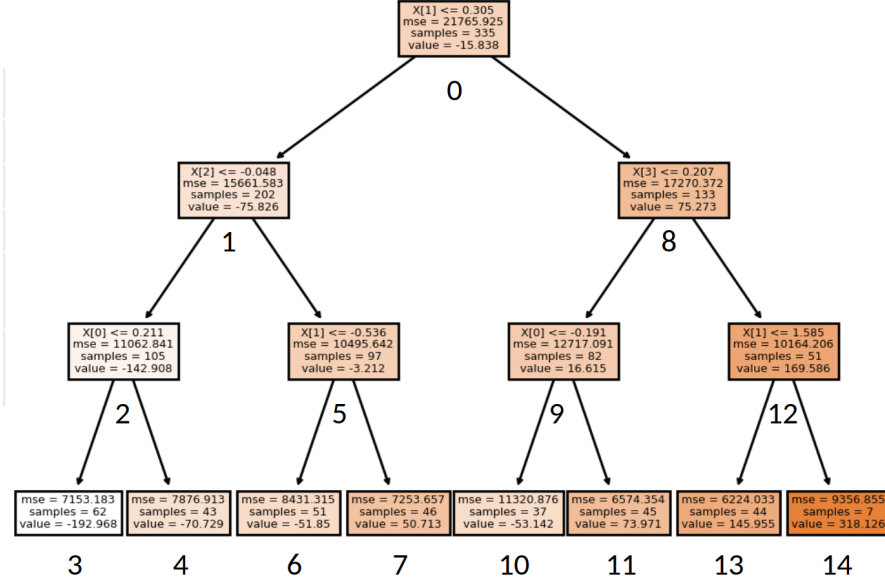


Figure 1: A simple decision tree used to illustrate the link between $\hat{f}^{(SHAP)}$ and Algorithm 1 in [6] (Tree SHAP) with numbered leaves.

Algorithm 1 Estimating $E[f(x) | x_S]$

```

procedure EXPVALUE( $x, S, tree = \{v, a, b, t, r, d\}$ )
  procedure G( $j, w$ )
    if  $v_j \neq \text{internal}$  then
      return  $w \cdot v_j$ 
    else
      if  $d_j \in S$  then
        return  $G(a_j, w)$  if  $x_{d_j} \leq t_j$  else  $G(b_j, w)$ 
      else
        return  $G(a_j, wr_{a_j}/r_j) + G(b_j, wr_{b_j}/r_j)$ 
      end if
    end if
  end procedure
  return  $G(1, 1)$ 
end procedure

```

Figure 2: Algorithm 1 in [6] (Tree SHAP)

$$\begin{aligned}
\hat{f}_{\{0,2\}}^{(SHAP)}(x_{ref}) &= P(X_1 \leq 0.305)P(X_2 > -0.048|X_1 \leq 0.305) * P(X_1 \leq -0.536|X_2 > -0.048)f_6 \\
&\quad + P(X_1 \leq 0.305)P(X_2 > -0.048|X_1 \leq 0.305) * P(X_1 > -0.536|X_2 > -0.048)f_7 \\
&\quad + P(X_1 > 0.305)P(X_3 \leq 0.207|X_1 > 0.305) * P(X_0 > -0.191|X_3 \leq 0.207)f_{11} \\
&\quad + P(X_1 > 0.305)P(X_3 > 0.207|X_1 > 0.305) * P(X_1 \leq 1.585|X_3 > 0.207)f_{13} \\
&\quad + P(X_1 > 0.305)P(X_3 > 0.207|X_1 > 0.305) * P(X_1 > 1.585|X_3 > 0.207)f_{14} \\
&= (202/335) * 1 * (51/97) * (-51.85) + (202/335) * 1 * (46/97) * (50.716) \\
&\quad + (133/335) * (82/133) * 1 * (73.971) + (133/335) * (51/133) * (44/51) * (145.955) \\
&\quad + (133/335) * (51/133) * (7/51) * (318.125) \\
&= 41.98
\end{aligned}$$

Step	Calculus
0	G(0, 1)
1	G(1, 202/335) + G(8, 133/335)
2	G(5, 202/335) + G(9, 88/335) + G(12, 51/335)
3	G(6, (202/335)*(51/97)) + G(7, (202/335)*(46/97)) + G(11, 82/335) + G(13, 44/335) + G(14, 7/335)
4	-(202/335)*(51/97)*51,85 + (202/335)*(46/97)*50,713 + (82/335)*73,971 + (44/335)*145,955 + (7/335)*318,126
5	= 41.98

Table 1: Step of Algorithm 1 in [6] (TreeSHAP) for the computation of $E[f(\mathbf{X}) | X_0 = 2, X_2 = 0.5]$ with $x_{ref} = [2, 3, 0.5, -1]$ and the tree in figure 1.

E Focus on influential variables in a linear model

Proposition E.1. *Let assumes that we have $X \in \mathbb{R}^p$, $X \in \mathcal{N}(0, I)$ and a linear predictor f defined as:*

$$f(X) = (a_1 X_1 + a_2 X_2) \mathbb{1}_{X_5 \leq 0} + (a_3 X_3 + a_4 X_4) \mathbb{1}_{X_5 > 0}. \quad (\text{E.1})$$

Even if we choose an observation \mathbf{x} such that $x_5 \leq 0$ and the predictor only uses X_1, X_2 , the SV of ϕ_3, ϕ_4 is not necessarily zero.

Proof.

$$\phi_3 = \frac{1}{p} \sum_{S \subseteq [p] \setminus \{3\}} \binom{p-1}{|S|}^{-1} \left(f_{S \cup 3}(\mathbf{x}_{S \cup 3}) - f_S(\mathbf{x}_S) \right) \quad (\text{E.2})$$

$$= \frac{1}{p} \sum_{S \subseteq [p] \setminus \{3,5\}} \binom{p-1}{|S|}^{-1} \left(f_{S \cup 3}(\mathbf{x}_{S \cup 3}) - f_S(\mathbf{x}_S) \right) + \frac{1}{p} \sum_{S \subseteq [p] \setminus \{3,5\}} \binom{p-1}{|S|+1}^{-1} \left(f_{S \cup \{3,5\}}(\mathbf{x}_{S \cup \{3,5\}}) - f_{S \cup 5}(\mathbf{x}_{S \cup 5}) \right) \quad (\text{E.3})$$

The second term is zero. Indeed, $\forall S \subseteq [p] \setminus \{3, 5\}$

$$f_{S \cup \{3,5\}}(\mathbf{x}_{S \cup \{3,5\}}) - f_{S \cup 5}(\mathbf{x}_{S \cup 5}) = 0.$$

Because, if we condition on the event $\{X_5 = \mathbf{x}_5\}$ with $x_5 \leq 0$

$$\begin{aligned} f_{S \cup \{3,5\}}(\mathbf{x}_{S \cup \{3,5\}}) &= E \left[(a_1 X_1 + a_2 X_2) \mathbb{1}_{X_5 \leq 0} + (a_3 X_3 + a_4 X_4) \mathbb{1}_{X_5 > 0} \mid X_{S \cup \{3,5\}} = \mathbf{x}_{S \cup \{3,5\}} \right] \\ &= E \left[(a_1 X_1 + a_2 X_2) \mathbb{1}_{X_5 \leq 0} \mid X_{S \cup \{3,5\}} = \mathbf{x}_{S \cup \{3,5\}} \right] \\ &= E \left[(a_1 X_1 + a_2 X_2) \mid X_{S \cup 5} = \mathbf{x}_{S \cup 5} \right] \\ &= f_{S \cup 5}(\mathbf{x}_{S \cup 5}). \end{aligned}$$

because $x_5 \leq 0$
independent of X_3

□

The first term of 3.3 is the classic marginal contribution of SV in linear model. $\forall S \subseteq [p] \setminus \{3, 5\}$,

$$\begin{aligned} f_{S \cup 3}(\mathbf{x}_{S \cup 3}) &= E \left[a_1 X_1 + a_2 X_2 \mid X_{S \cup 3} = \mathbf{x}_{S \cup 3} \right] P(X_5 \leq 0 \mid X_{S \cup 3} = \mathbf{x}_{S \cup 3}) \\ &\quad + E \left[a_3 X_3 + a_4 X_4 \mid X_{S \cup 3} = \mathbf{x}_{S \cup 3} \right] P(X_5 > 0 \mid X_{S \cup 3} = \mathbf{x}_{S \cup 3}) \\ &= E \left[a_1 X_1 + a_2 X_2 \mid X_S = \mathbf{x}_S \right] P(X_5 \leq 0) + (E \left[a_2 X_2 \mid X_S = \mathbf{x}_S \right] + a_3 \mathbf{x}_3) P(X_5 > 0) \\ &= f_S(\mathbf{x}_S) + P(X_5 > 0) \left(a_3 (\mathbf{x}_3 - E[X_3]) \right) \end{aligned}$$

Therefore,

$$\begin{aligned} \phi_3 &= \frac{1}{p} \sum_{S \subseteq [p] \setminus \{3,5\}} \binom{p-1}{|S|}^{-1} P(X_5 > 0) \left(a_3 (\mathbf{x}_3 - E[X_3]) \right) \\ &= K \left(a_3 (\mathbf{x}_3 - E[X_3]) \right) \end{aligned} \quad K \text{ is a constant}$$

The computation of ϕ_4 is symmetric.

F Additional examples

F.1 Impact of quantile discretization

The table F.2 (borrowed from [1]) shows the impact of discretization on the performance of a Random Forest on UCI datasets.

Dataset	Breiman's RF	q=2	q=5	q=10	q=20
Authentication	0.0002	0.08	0.002	0.0005	0.0004
Diabetes	0.17	0.23	0.18	0.18	0.18
Haberman	0.32	0.35	0.30	0.32	0.30
Heart Statlog	0.10	0.10	0.10	0.10	0.10
Hepatitis	0.13	0.15	0.14	0.14	0.13
Ionosphere	0.02	0.07	0.03	0.02	0.02
Liver Disorders	0.23	0.32	0.27	0.25	0.24
Sonar	0.07	0.09	0.07	0.07	0.07
Spambase	0.01	0.14	0.03	0.02	0.01
Titanic	0.13	0.15	0.14	0.14	0.13
Wilt	0.007	0.15	0.03	0.02	0.02

Table 2: Accuracy, measured by 1-AUC on UCI datasets, for two algorithms: Breiman's random forests and random forests with splits limited to q-quantiles, for $q \in \{2, 5, 10, 20\}$. Table 5 in [1]

In the table F.3 below, we compare the different estimators (Tree SHAP, Leaf and Discrete) after the discretization of continuous variables. The model and the data used are described in Section A.2. We see that the Discrete estimator is much more sensitive to out-of-distribution than Leaf. It can also be observed that the different estimators have considerable variance. Indeed, a model trained on discrete variables tends to have poorly filled leaves.

Datasets	$D^{(Explain)}$		Uniform	
Metrics	MSE	TPR	MSE	TPR
Tree SHAP	2.92	55% (49%)	2.91	74% (43%)
Leaf	0.98	74% (43%)	2.64	75% (43%)
Discrete	0.20	95% (22%)	2.99	65.79% (47%)

Table 3: Metrics of the different estimators after discretization

F.2 The differences between Coalition and sum on Census Data

We use UCI Adult Census Dataset [4]. We keep only 4 highly-predictive categorical variables: Marital Status, Workclass, Race, Education and use a Random Forest which has a test accuracy of 86%. We compare the Global SV by taking the coalition or sum of the modalities over N=5000 observations. Global SV are defined as:

$$I_j = \sum_{i=0}^N |\phi_j^{(i)}|$$

In figure 3, we see differences between the global SV with coalition and sum. The ranking of the variables changes, e.g. Education goes from important with sum to not important with the coalition.

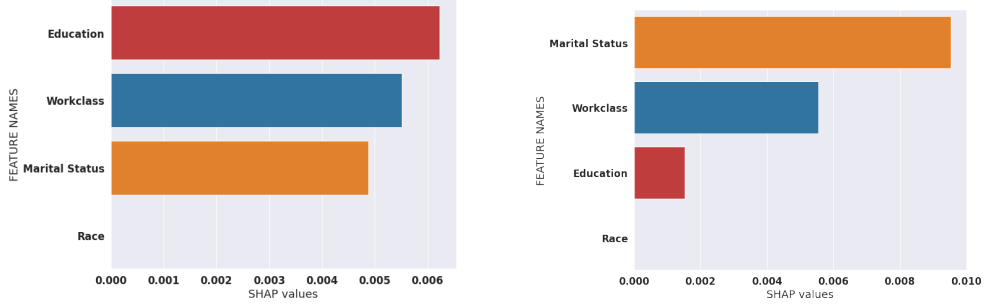


Figure 3: Difference between the global absolute value of SV: sum (left) vs coalition (right) of dummies

We also compute the proportion of order inversion over $N=5000$ observations choose randomly. The ranking of variables is changed in 10% of the cases. Note that this difference may increase or diminish depending on the data.

F.3 SDP and Active SV analysis on Lucas Data

We use an accurate decision tree trained on LUCAS [5], a dataset generated by causal Bayesian networks with 12 binary variables. The graph is drawn in figure 4 and we provide the probability table in figure 5

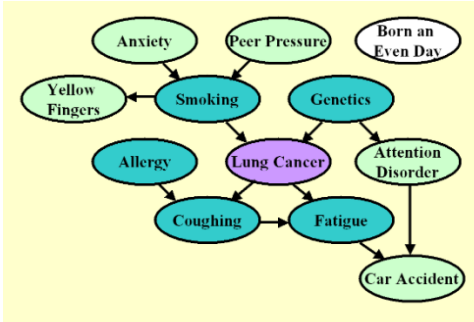


Figure 4: Bayesian network that represents the causal relationships between variables

$P(\text{Anxiety}=\text{T})=0.64277$
$P(\text{Peer Pressure}=\text{T})=0.32997$
$P(\text{Smoking}=\text{T} \text{Peer Pressure}=\text{F}, \text{Anxiety}=\text{F})=0.43118$
$P(\text{Smoking}=\text{T} \text{Peer Pressure}=\text{T}, \text{Anxiety}=\text{F})=0.74591$
$P(\text{Smoking}=\text{T} \text{Peer Pressure}=\text{F}, \text{Anxiety}=\text{T})=0.8686$
$P(\text{Smoking}=\text{T} \text{Peer Pressure}=\text{T}, \text{Anxiety}=\text{T})=0.91576$
$P(\text{Yellow Fingers}=\text{T} \text{Smoking}=\text{F})=0.23119$
$P(\text{Yellow Fingers}=\text{T} \text{Smoking}=\text{T})=0.95372$
$P(\text{Genetics}=\text{T})=0.15953$
$P(\text{Lung cancer}=\text{T} \text{Genetics}=\text{F}, \text{Smoking}=\text{F})=0.23146$
$P(\text{Lung cancer}=\text{T} \text{Genetics}=\text{T}, \text{Smoking}=\text{F})=0.86996$
$P(\text{Lung cancer}=\text{T} \text{Genetics}=\text{F}, \text{Smoking}=\text{T})=0.83934$
$P(\text{Lung cancer}=\text{T} \text{Genetics}=\text{T}, \text{Smoking}=\text{T})=0.99351$
$P(\text{Attention Disorder}=\text{T} \text{Genetics}=\text{F})=0.28956$
$P(\text{Attention Disorder}=\text{T} \text{Genetics}=\text{T})=0.68706$
$P(\text{Born an Even Day}=\text{T})=0.5$
$P(\text{Allergy}=\text{T})=0.32841$
$P(\text{Coughing}=\text{T} \text{Allergy}=\text{F}, \text{Lung cancer}=\text{F})=0.1347$
$P(\text{Coughing}=\text{T} \text{Allergy}=\text{T}, \text{Lung cancer}=\text{F})=0.64592$
$P(\text{Coughing}=\text{T} \text{Allergy}=\text{F}, \text{Lung cancer}=\text{T})=0.7664$
$P(\text{Coughing}=\text{T} \text{Allergy}=\text{T}, \text{Lung cancer}=\text{T})=0.99947$
$P(\text{Fatigue}=\text{T} \text{Lung cancer}=\text{F}, \text{Coughing}=\text{F})=0.35212$
$P(\text{Fatigue}=\text{T} \text{Lung cancer}=\text{T}, \text{Coughing}=\text{F})=0.56514$
$P(\text{Fatigue}=\text{T} \text{Lung cancer}=\text{F}, \text{Coughing}=\text{T})=0.80016$
$P(\text{Fatigue}=\text{T} \text{Lung cancer}=\text{T}, \text{Coughing}=\text{T})=0.89589$
$P(\text{Car Accident}=\text{T} \text{Attention Disorder}=\text{F}, \text{Fatigue}=\text{F})=0.2274$
$P(\text{Car Accident}=\text{T} \text{Attention Disorder}=\text{T}, \text{Fatigue}=\text{F})=0.779$
$P(\text{Car Accident}=\text{T} \text{Attention Disorder}=\text{F}, \text{Fatigue}=\text{T})=0.78861$
$P(\text{Car Accident}=\text{T} \text{Attention Disorder}=\text{T}, \text{Fatigue}=\text{T})=0.97169$

Figure 5: Probabilities table used to generate Data

We want to explain an observation with a well-defined ground truth. We know from the probability tables that if Smoking, Genetic, Coughing are False, the probability of having Cancer is very low. So, we should have these three variables in the Sufficient Coalition: this is what we can observe in table F.4.

We have also computed the Active SV and the standard SV. The figure 6 shows that the Active SV are indeed sparse giving importance to the local active SV while standard SV found that Fatigue, Yellow Fingers, Anxiety is more important of Genetic for this observation.

Active and Null coalition	SDP
$S_\pi^*(x) = [\text{Smoking, Genetics, Coughing}]$	0.96
$N_\pi(x) = [\text{Yellow Fingers, Anxiety, Peer Pressure, Attention Disorder, Born an Even Day, Car Accident, Fatigue, Allergy}]$	0.77

Table 4: The Sufficient coalition found with $\pi = 0.9$

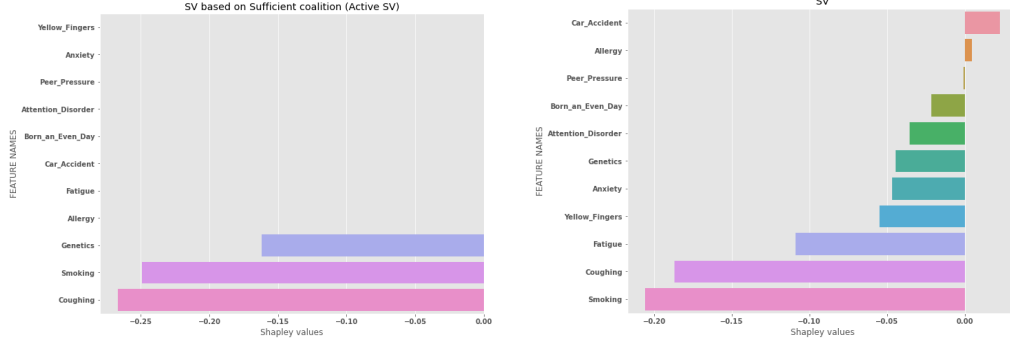


Figure 6: Left figure: SV ϕ_i^* computed with the Sufficient Coalition given in figure 6. Right figure: SV ϕ_i computed with all the variables.

G Individual Shapley values for indicator variables in Dummy Encoding

We give some partial results for the Shapley Values of the modalities $Y = k$, based on the dummy encoding considered in section 2. Indeed equation 2.4 introduces $\phi_k(\tilde{f}, x, y_{1:K-1})$, and proposition 2.1 claims that their sum is different in all generality of the SV of Y . In this section, we give a deeper insight into these values and show that are related multiple comparisons between modalities.

We compute the Shapley Value at point $(x, y = i) = (x, 0, 0, \dots, 1, \dots, 0) = (x, \mathcal{C}(y))$: for ease of notation, we set $Y_0 = X$, and we compute also the Shapley values $\phi_k(\tilde{f}; x, y_{1:K-1})$ for $k = 1, \dots, K-1$. We recall that we need to compute

$$\frac{1}{K} \sum_{k=0}^{K-1} \frac{1}{\binom{K-1}{k}} \sum_{\substack{Z \subseteq \llbracket 1..K \rrbracket \setminus i \\ |Z| = k}} \Delta(\tilde{f}; Z, i).$$

where Δ denotes the difference between the value function evaluated at $Z \cup \{i\}$ and Z . If we examine the terms $\Delta(\tilde{f}; Z, i)$, the computation needs to take into account if $X = Y_0$ is part of the conditioning variable of not. Indeed, we have for each $k \geq 1$,

$$\sum_{\substack{Z \subseteq \llbracket 0..K-1 \rrbracket \setminus i \\ |Z| = k}} \Delta(\tilde{f}; Z, i) = \sum_{\substack{Z \subseteq \llbracket 1..K-1 \rrbracket \setminus i \\ |Z| = k}} \Delta(\tilde{f}; Z, i) + \sum_{\substack{Z' \subseteq \llbracket 1..K-1 \rrbracket \setminus i \\ |Z'| = k-1}} \Delta(\tilde{f}; Z' \cup \{0\}, i). \quad (\text{G.1})$$

We start by computing the first term in the right hand side, and it involves only the dummies, and not the quantitative variable.

Proposition G.1 (Computation of Contributions in Shapley without X). *We compute the Shapley values of the variable Y_i , when we have the observations $(x, y_{1:K-1}) = (x, \mathcal{C}(i))$ for $i \in \{1, \dots, K\}$. We consider any $Z' \subseteq \llbracket 1..K-1 \rrbracket \setminus i$, with $|Z'| = k \geq 1$ and $Z' = \{j_1, \dots, j_k\}$. In that case,*

$$\Delta(\tilde{f}; Z, i) = E_P[f(X, Y)|Y = i] - E_P[f(X, Y)|Y \notin \{j_1, \dots, j_k\}] \quad (\text{G.2})$$

Proof. We have $Y_i = 1 \Leftrightarrow Y = i$, and for $Z' \subseteq \llbracket 1..K-1 \rrbracket \setminus \{0, i\}$, we consider $Z' = \{j_1, \dots, j_k\}$, with $1 \leq j_1 < \dots < j_k \leq K-1$,

$$\begin{aligned} E_{\tilde{P}} \left[\tilde{f}(Y_0, Y_{1:K-1}) | Y_{j_1} = 0, \dots, Y_{j_k} = 0, Y_i = 1 \right] &= E_{\tilde{P}} \left[\tilde{f}(Y_0, Y_{1:K-1}) | Y_i = 1 \right] \\ &= E_{\tilde{P}} \left[\tilde{f}(Y_0, \mathcal{C}(Y)) | Y_i = 1 \right] \\ &= E_P [f(Y_0, Y) | Y = i] \end{aligned}$$

because for all $j_1, \dots, j_{k-1} \neq i$, we have $\{Y_{j_1} = 0, \dots, Y_{j_k} = 0, Y_i = 1\} = \{Y_i = 1\}$.
Moreover,

$$E_{\tilde{P}} \left[\tilde{f}(Y_0, Y_{1:K-1}) | Y_{j_1} = 0, \dots, Y_{j_k} = 0 \right] = E_P \left[\tilde{f}(Y_0, \mathcal{C}(Y)) | Y \neq j_1, \dots, j_k \right]$$

Hence for $Z \subseteq \llbracket 1..K-1 \rrbracket \setminus i$, we have

$$\Delta(\tilde{f}; Z, i) = E_P [f(X, Y) | Y = i] - E_P [f(X, Y) | Y \notin \{j_1, \dots, j_k\}].$$

□

The second term of the right hand side is given below.

Proposition G.2 (Computation of Contributions in Shapley with X). *We compute the Shapley values only for the variable Y_i , when we have the observations doable $(x, y_{1:K-1}) = (x, \mathcal{C}(i))$ for $i \in \{1, \dots, K\}$. We consider any $Z' \subseteq \llbracket 1..K-1 \rrbracket \setminus i$, with $|Z'| = k-1 \geq 1$, and $Z' = \{j_1, \dots, j_{k-1}\}$. In that case,*

$$\Delta(\tilde{f}; Z' \cup \{0\}, i) = E_P [f(X, Y) | X = x, Y = i] - E_P [f(X, Y) | X, Y \notin \{j_1, \dots, j_{k-1}\}] \quad (\text{G.3})$$

Proof. We assume that we have a subset $|Z'| = k-1$, such that $Z' \subseteq \llbracket 1..K-1 \rrbracket \setminus i$. This means that $Z' = \{j_1, \dots, j_{k-1}\}$, with $1 \leq j_1, \dots, j_{k-1} \leq K-1$. We

$$\begin{aligned} E_{\tilde{P}} \left[\tilde{f}(Y_0, Y_{1:K-1}) | Y_0 = x, Y_{j_1} = 0, \dots, Y_{j_{k-1}} = 0, Y_i = 1 \right] &= E_{\tilde{P}} \left[\tilde{f}(Y_0, Y_{1:K-1}) | Y_0 = x, Y_i = 1 \right] \\ &= E_P \left[\tilde{f}(Y_0, \mathcal{C}(Y)) | Y_0 = x, Y = i \right] \\ &= E_P [f(Y_0, \mathcal{C}(Y)) | Y_0 = x, Y = i] \end{aligned}$$

and

$$\begin{aligned} E_{\tilde{P}} \left[\tilde{f}(Y_0, Y_{1:K-1}) | Y_0 = x, Y_{j_1} = 0, \dots, Y_{j_{k-1}} = 0 \right] &= E_P \left[\tilde{f}(Y_0, \mathcal{C}(Y)) | Y_0 = x, Y \notin \{j_1, \dots, j_{k-1}\} \right] \\ &= E_P [f(Y_0, Y) | Y_0 = x, Y \notin \{j_1, \dots, j_{k-1}\}] \end{aligned}$$

□

Finally, we can give several examples of the different increments involved in the Shapley values of each variable X or Y_k . If $k = 0$, then $Z' = \emptyset$ and

$$\Delta(\tilde{f}; Z', i) = \Delta(\tilde{f}; \emptyset, i) = E_P [f(X, Y) | Y = i] - E_P [f(X, Y)]$$

If $k = 1$, then $Z' = \{0\}$ or $Z' = \{j\} \neq \{i\}$,

$$\Delta(\tilde{f}; Z', i) = \Delta(\tilde{f}; 0, i) = E_P [f(X, Y) | X = x, Y = i] - E_P [f(X, Y) | X = x]$$

$$\Delta(\tilde{f}; Z', i) = \Delta(\tilde{f}; \{j\}, i) = E_P [f(X, Y) | Y = i] - E_P [f(X, Y) | Y \neq j]$$

For $k = K-1$, $Z' = \{1, \dots, K-1\}$,

$$\Delta(\tilde{f}; \{1, \dots, K-1\}, i) = E_P [f(X, Y) | X = x, Y = i] - E_P [f(X, Y) | X = x, Y \neq i]$$

The propositions G.1 and G.2 show that the individual Shapley value for the variable (modality) Y_i is a weighted mean of the difference between classe i and group of classes:

$$\begin{cases} E_P [f(X, Y) | Y = i] - E_P [f(X, Y) | Y \notin \{j_1, \dots, j_k\}] \\ E_P [f(X, Y) | X = x, Y = i] - E_P [f(X, Y) | X, Y \notin \{j_1, \dots, j_{k-1}\}] \end{cases}$$

Finally, we can also compute the Shapley values of the other variables Y_j at point $(x, y = i)$, for $j \neq i$. In that case, the difference $\Delta(\tilde{f}; Z', j), j \neq i$ are of the type of

$$\begin{cases} E_P[f(X, Y)|Y \notin \{j, j_1, \dots, j_k\}] - E_P[f(X, Y)|Y \notin \{j_1, \dots, j_k\}] \\ E_P[f(X, Y)|Y = i] - E_P[f(X, Y)|Y = i] \\ E_P[f(X, Y)|X = x, Y \notin \{j, j_1, \dots, j_k\}] - E_P[f(X, Y)|X, Y \notin \{j_1, \dots, j_{k-1}\}] \\ E_P[f(X, Y)|X = x, Y = i] - E_P[f(X, Y)|X, Y = i] \end{cases}$$

The Shapley values computes a mean of the difference between different aggregation of modalities, that contains or not the variable of interest.

As a conclusion of this part, we see that the individual Shapley values $\phi_k(\tilde{f}; x, y_{1:K-1})$ perform a multiple comparison of the means obtained by aggregating the classes or modalities in various ways, looking at the presence or not of the modality k . These differences of means have weights $\frac{1}{\binom{K-1}{k}}$ where k is basically the number of classes of the variable Y that we aggregate.

Consequently the sum $\sum_{k=1}^K \phi_k(\tilde{f}; x, y_{1:K-1})$ is clearly different from the

$$\phi_Y(f; x, y) = \frac{1}{2} (E[f(X, Y)|Y = y] - E[f(X, Y)]) + \frac{1}{2} (f(x, y) - E[f(X, Y)|X = x]).$$

This latter has a much more global analysis that aims at measuring how the mean $E[f(X, Y)|Y = y]$ in the various classes changes w.r.t $E[f(X, Y)]$, while the individual Shapley focus on the difference between subgroups of classes.

H Plug-In estimator of Marginal expectation

As we have indicated in the paper, the Shapley Values can be computed with different probability $Q_{S, \mathbf{x}}$. In that section, we show that when we use the marginal distribution (as in the so-called interventional case), the previous estimators for tree-based models can be adapted straightforwardly. We consider then decision tree

$$f(x) = \sum_{m=1}^M f_m \mathbb{1}_{L_m}(x)$$

and remark that the Marginal Shapley coefficients involve the computations of the marginal expectations $E_P[\mathbb{1}_{L_m}(\mathbf{X}_{\bar{Z}}, \mathbf{x}_Z)]$ for any subgroup of variables Z . On real data, we need to compute the conditional expectations, but we use the Tree approximations in order to replace

$$\begin{aligned} E_P[\mathbb{1}_{L_m}(\mathbf{X}_{\bar{Z}}, \mathbf{x}_Z)] &= \int \int \mathbb{1}_{L_m}(\mathbf{u}_{\bar{Z}}, \mathbf{x}_Z) p(\mathbf{u}_{\bar{Z}}, \mathbf{u}_Z) d\mathbf{u}_{\bar{Z}} d\mathbf{u}_Z \\ &= \int \int \mathbb{1}_{L_m}(\mathbf{u}_{\bar{Z}}, \mathbf{x}_Z) p(\mathbf{u}_Z | \mathbf{u}_{\bar{Z}}) p(\mathbf{u}_{\bar{Z}}) d\mathbf{u}_Z d\mathbf{u}_{\bar{Z}} \\ &= \int \left\{ \int p(\mathbf{u}_Z | \mathbf{u}_{\bar{Z}}) d\mathbf{u}_Z \right\} \mathbb{1}_{L_m}(\mathbf{u}_{\bar{Z}}, \mathbf{x}_Z) p(\mathbf{u}_{\bar{Z}}) d\mathbf{u}_{\bar{Z}} \\ &= \int \mathbb{1}_{L_m}(\mathbf{u}_{\bar{Z}}, \mathbf{x}_Z) p(\mathbf{u}_{\bar{Z}}) d\mathbf{u}_{\bar{Z}} \end{aligned}$$

This means that we just need the marginal distributions of the variables $\mathbf{X}_{\bar{Z}}$ in order to compute the expectations of the leaf. In the case of quantitative data, the leaf can be written $L_m = \prod_{i=1}^p [a_i^m, b_i^m]$, and we have by definition

$$\exists k \in Z, x_k \notin [a_k, b_k] \implies \mathbb{1}_{L_m}(\mathbf{u}_{\bar{Z}}, \mathbf{x}_Z) = 0$$

We define the set of leafs compatible with condition $\mathbf{X}_Z = \mathbf{x}_Z$ as

$$C(Z, \mathbf{x}) = \left\{ m \in [1 \dots M] \mid L_m = \prod_{i=1}^p [a_i^m, b_i^m], \forall k \in Z, x_k \in [a_k^m, b_k^m] \right\}$$

We write for $m \in C(Z, \mathbf{x})$, $L_m = L_m^{\bar{Z}} \times L_m^Z$, with $L_m^{\bar{Z}} = \prod_{i \in \bar{Z}} [a_i^m, b_i^m]$ and $L_m^Z = \prod_{i \in Z} [a_i^m, b_i^m]$, this means that for all $m \in C(Z, \mathbf{x})$ we have

$$E_P [\mathbb{1}_{L_m}(\mathbf{X}_{\bar{Z}}, \mathbf{x}_Z)] = E_P [\mathbb{1}_{L_m^{\bar{Z}}}(\mathbf{X}_{\bar{Z}})]$$

As an approximation, the conditional probability for $m \in C(Z, \mathbf{x})$ is computed as

$$\begin{aligned} E_P [\mathbb{1}_{L_m^{\bar{Z}}}(\mathbf{X}_{\bar{Z}})] &= P(X_i \in [a_i^m, b_i^m], i \in \bar{Z}) \\ &\simeq \frac{N(L_m^{\bar{Z}})}{N} \end{aligned}$$

where $N(L_m^{\bar{Z}})$ is the number of observations in the (partial) leaf $L_m^{\bar{Z}}$. As a consequence we have

$$\begin{aligned} E_P [f(\mathbf{X}_{\bar{Z}}, \mathbf{x}_Z)] &= \sum_{m=1}^M \hat{y}_m E_P [\mathbb{1}_{L_m}(\mathbf{X}_{\bar{Z}}, \mathbf{x}_Z)] \\ &= \sum_{m \in C(Z, \mathbf{x})} \hat{y}_m E_P [\mathbb{1}_{L_m}(\mathbf{X}_{\bar{Z}}, \mathbf{x}_Z)] \\ &= \sum_{m \in C(Z, \mathbf{x})} \hat{y}_m E_P [\mathbb{1}_{L_m^{\bar{Z}}}(\mathbf{X}_{\bar{Z}})] \\ &\simeq \sum_{m \in C(Z, \mathbf{x})} \hat{y}_m \frac{N(L_m^{\bar{Z}})}{N} \end{aligned}$$

I EXPERIMENTAL SETTINGS

All our experiments are reproducible and can be found on the github repository *Active Coalition of Variables*, www.github.com/salimamoukou/acv00

A.1 Toy model of Section 2.3

Recall that the model is a linear predictor with categorical variables define as $f(X, Y) = B_Y X$ with $X|Y = y \sim \mathcal{N}(\mu_y, \Sigma_y)$ and $\mathbb{P}(Y = y) = \pi_y$, $Y \in \{a, b, c\}$.

For the experiments in Figure 1 and 2, we set $\pi_y = \frac{1}{3}$, $\mu_y = 0 \forall y \in \{a, b, c\}$. We use a random matrices generated from a Wishart distribution. The covariance matrices are:

$$\begin{aligned} \Sigma_a &= \begin{bmatrix} 0.41871254 & -0.790061361 & 0.46956991 \\ -0.79006136 & 1.90865098 & -0.82571655 \\ 0.46956991 & -0.82571655 & 0.95835472 \end{bmatrix}, \Sigma_b = \begin{bmatrix} 0.55326081 & 0.11811951 & -0.70677924 \\ 0.11811951 & 2.73312979 & -2.94400196 \\ -0.70677924 & -2.94400196 & 4.22105088 \end{bmatrix}, \\ \Sigma_c &= \begin{bmatrix} 9.2859966 & 1.12872646 & 2.4224434 \\ 1.12872646 & 0.92891237 & -0.14373393 \\ 2.4224434 & -0.14373393 & 1.81601676 \end{bmatrix} \text{ for } y \in \{a, b, c\} \text{ respectively.} \end{aligned}$$

The coefficients are $B_a = [1, 3, 5]$, $B_b = [-5, -10, -8]$, $B_c = [6, 1, 0]$ and the selected observation in figure 1 values is $x = [0.35, -1.61, -0.11, 1., 0., 0.]$

A.2 Toy model of Section 3.3

The data $\mathcal{D}_x^{(Train)} = (x_i, z_i)_{1 \leq i \leq n}$ are generated from a linear regression $Z = B^t X$ with $n = 10000$, $X \sim \mathcal{N}(\mu, \Sigma)$ where $\mu = 0$, $\Sigma = 0.7 \times \text{np.ones}(d, d) - (0.7-1) \times \text{np.eyes}(7)$. $d=5$, $B = [6.49, -2.44, -2.11, -4.29, 3.46]$ for the continuous case and $d=3$, $B = [6.49, -2.44, 0]$ for the discrete case.

We used a decision tree on $\mathcal{D}^{(Train)}$ with the defaults parameters. The Mean Squared Error (MSE) are $\text{MSE} = 4.39$ for the continuous case and $\text{MSE} = 2.88$ for the discrete case.

A.3 Comparisons of calculation time of SV estimates with ACV and TreeSHAP

We show below a run-time comparison of the computation of n SV with ACV and TreeSHAP. We used 3 datasets with different shape: Boston (N=506, p=13), Adults (N=32561, p=12), Toy linear

model ($N=50000$, $p=500$). The model used was XGBoost with default parameters ($ntree = 100$, $maxdepth = 6$). We compute the SV of $n=1000$ observations for Adults, Toy model and $n=506$ for Boston.

Datasets	Boston($n = 506, p = 13$)	Adults ($n = 1000, p = 12$)
Leaf	8.82 s (204 ms)	1 min 4 s (1.73 s)
Tree SHAP	129 ms (6.91 ms)	3.33 s (39.9 ms)

This difference in runtime can be partly explained by the fact that Leaf estimator has to go through all the data for each leaf, whereas TreeSHAP uses the information stored in the tree.

References

- [1] Clément Bénard, Gérard Biau, Sébastien Da Veiga, and Erwan Scornet. Sirius: Stable and interpretable rule set for classification. *Electronic Journal of Statistics*, 15(1):427–505, 2021.
- [2] S. Chen, Arthur Choi, and Adnan Darwiche. The same-decision probability: A new tool for decision making, 2012.
- [3] Suming Chen, Arthur Choi, and Adnan Darwiche. An exact algorithm for computing the same-decision probability. *IJCAI ’13*, page 2525–2531. AAAI Press, 2013.
- [4] Dheeru Dua and Casey Graff. UCI machine learning repository, 2017.
- [5] Dataset LUCAS. Lucas (lung cancer simple set) dataset. <http://www.causality.inf.ethz.ch/data/LUCAS.html>.
- [6] Scott M. Lundberg, Gabriel Erion, Hugh Chen, Alex DeGrave, Jordan M. Prutkin, Bala Nair, Ronit Katz, Jonathan Himmelfarb, Nisha Bansal, and Su-In Lee. From local explanations to global understanding with explainable ai for trees. *Nature Machine Intelligence*, 2(1):2522–5839, 2020.
- [7] Eric Wang, Pasha Khosravi, and Guy Van den Broeck. Towards probabilistic sufficient explanations. In *Extending Explainable AI Beyond Deep Models and Classifiers Workshop at ICML (XXAI)*, 2020.