

Beyond TreeSHAP: Efficient Computation of Any-Order Shapley Interactions for Tree Ensembles

Maximilian Muschalik^{1,*}, Fabian Fumagalli^{2,*}, Barbara Hammer², Eyke Hüllermeier¹

¹LMU Munich, MCML Munich, D-80539 Munich, Germany

²Bielefeld University, CITEC, D-33619 Bielefeld, Germany

maximilian.muschalik@lmu.de, ffumagalli@techfak.uni-bielefeld.de, bhammer@techfak.uni-bielefeld.de, eyke@lmu.de

Abstract

While shallow decision trees may be interpretable, larger ensemble models like gradient-boosted trees, which often set the state of the art in machine learning problems involving tabular data, still remain black box models. As a remedy, the Shapley value (SV) is a well-known concept in explainable artificial intelligence (XAI) research for quantifying additive feature attributions of predictions. The model-specific TreeSHAP methodology solves the exponential complexity for retrieving exact SVs from tree-based models. Expanding beyond individual feature attribution, Shapley interactions reveal the impact of intricate feature interactions of any order. In this work, we present TreeSHAP-IQ, an efficient method to compute any-order additive Shapley interactions for predictions of tree-based models. TreeSHAP-IQ is supported by a mathematical framework that exploits polynomial arithmetic to compute the interaction scores in a single recursive traversal of the tree, akin to Linear TreeSHAP. We apply TreeSHAP-IQ on state-of-the-art tree ensembles and explore interactions on well-established benchmark datasets.

1 Introduction

Tree-based ensemble methods, in particular gradient-boosted trees (Friedman 2001), such as XGBoost (Chen and Guestrin 2016) or LightGBM (Ke et al. 2017), are among the most popular machine learning (ML) models and often achieve state-of-the-art (SOTA) performance on tabular data without extensive hyperparameter tuning (Shwartz-Ziv and Armon 2022). These ensemble methods utilize intricate prediction functions by employing tree structures of high depth, thereby obstructing interpretation of the model’s internal reasoning. Yet, understanding a model’s prediction is necessary for safe and reliable deployment, alongside addressing ethical and regulatory considerations (Adadi and Berrada 2018). Additive feature attributions, which split the individual features’ contributions to the prediction, are a prevalent approach to improving the local interpretation of ML models (Lundberg and Lee 2017; Covert and Lee 2021; Chen et al. 2023). However, in complex real-world applications, such as bioinformatics (Lunetta et al. 2004; Boulesteix et al. 2012; Winham et al. 2012; Wright, Ziegler, and König 2016) or language-related tasks (Tsang, Rambhatla, and Liu 2020)

*These authors contributed equally.
Accepted at AAAI’24.

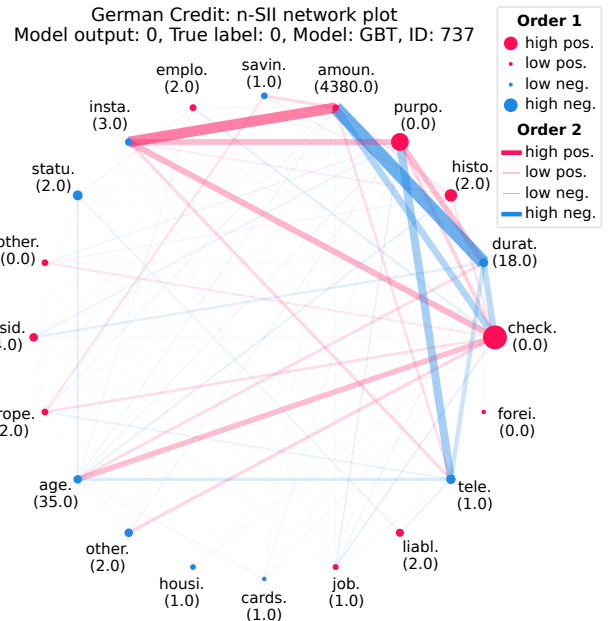


Figure 1: Network Plot after (Inglis, Parnell, and Hurley 2022) for a test instance of the *German Credit* dataset for visualizing local feature attribution and interaction.

features only attain meaningfulness when *interacting* with other features. In such scenarios, information about interactions complements additive feature attributions, which only show part of the picture (Wright, Ziegler, and König 2016).

In this work, we are interested in model-specific local XAI measures for tree-based models, such as XGBoost. In particular, the extension of predominant attribution measures based on the Shapley value (SV) (Shapley 1953) to any-order additive Shapley-based interactions to explain single predictions locally. Our work extends path dependent TreeSHAP (Lundberg et al. 2020), which exploits the structure of trees to reduce time complexity from exponential to polynomial, to any-order Shapley-based interactions.

Related Work. The SV (Shapley 1953) is a concept from cooperative game theory that has been proposed for model-agnostic explanations for local (Strumbelj and Kononenko 2014; Lundberg and Lee 2017) and global (Casalicchio,

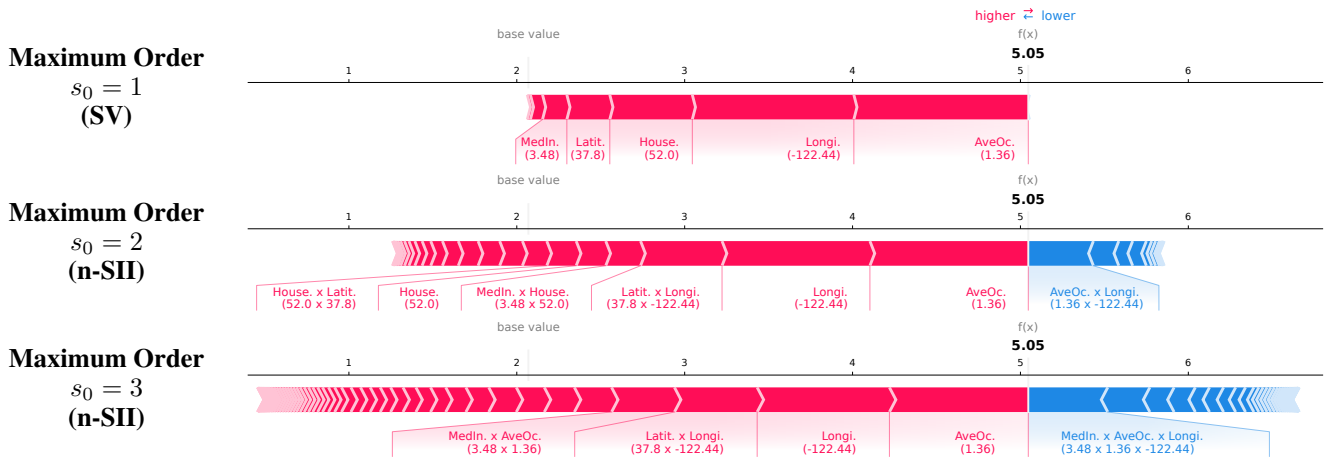


Figure 2: Force plots of positive (red) and negative (blue) SVs and n-SII scores for an instance of the *California* dataset The *longit.* feature has a high contribution, describing the proximity to the ocean, which affects the price. *TreeSHAP* ($s_0 = 1$) reveals this contribution. It also shows that *latit.* contributed positively. *TreeSHAP-IQ*, e.g. $s_0 \geq 2$, reveals that this contribution can be (mostly) attributed to the interaction *latit. x longit.*, which reveals that the *exact location*, and not *latit.*, is meaningful.

Molnar, and Bischl 2018; Covert, Lundberg, and Lee 2020) interpretation. In a model-agnostic setting, efficient approximations techniques, based on Monte Carlo (Castro, Gómez, and Tejada 2009; Castro et al. 2017; Kolpaczki et al. 2023; Fumagalli et al. 2023) or the representation of the SV as a constrained weighted least square problem (Lundberg and Lee 2017; Covert and Lee 2021; Jethani et al. 2022) have been proposed to overcome the exponential complexity. For tree-based models the SV can be computed in polynomial time using *TreeSHAP* (Lundberg et al. 2020) with more efficient variants (Yang 2021). Linear *TreeSHAP* (Yu et al. 2022) establishes a theoretical foundation that connects the computation to polynomial arithmetic, achieving SOTA computational and storage efficiency.

Limitations of the SV due to correlations and interactions have been widely studied by Slack et al. (2020), Sundararajan and Najmi (2020), and Kumar et al. (2020, 2021). Extensions to interactions have been proposed with the *Shapley Interaction Index* (SII) (Grabisch and Roubens 1999), its aggregation as *n-Shapley Values* (n-SII) (Bordt and von Luxburg 2023), the *Shapley Taylor Interaction Index* (STI) (Sundararajan, Dhamdhere, and Agarwal 2020) and the *Faithful Shapley Interaction Index* (FSI) (Tsai, Yeh, and Ravikumar 2023). All of these are subsumed in the broad class of the *Cardinal Interaction Index* (CII) (Grabisch and Roubens 1999). Model-agnostic approximations have been proposed for general CII (Fumagalli et al. 2023), STI (Sundararajan, Dhamdhere, and Agarwal 2020), SII and for FSI (Tsai, Yeh, and Ravikumar 2023). Local pairwise interactions for tree-based models were computed by Lundberg et al. (2020) and for interventional SHAP by Zern, Broelemann, and Kasneci (2023).

Other interaction scores were introduced by Tsang, Rambhatla, and Liu (2020), Zhang et al. (2021), Patel, Strobel, and Zick (2021), Harris, Pymar, and Rowat (2022), and Hiabu, Meyer, and Wright (2023). Interaction scores are

further linked to functional decomposition (Hooker 2004, 2007; Lengerich et al. 2020; Herbringer, Bischl, and Casalicchio 2023). For tree-based models, limitations of feature attribution measures (Wright, Ziegler, and König 2016), and efficient implementations for interactions (Lengerich et al. 2020; Hiabu, Meyer, and Wright 2023) were discussed.

So far, any-order Shapley interactions have only been studied in a model-agnostic setting, where the exponential complexity problem is approximately solved. Tree-based approaches have not considered the efficient computation of local any-order Shapley interactions.

Contribution. Our main contributions include;

1. *TreeSHAP-IQ* (Section 3): An efficient algorithm for computing any-order SII scores for tree ensembles. *TreeSHAP-IQ* is supported by a mathematical framework based on polynomial arithmetic, akin to Linear *TreeSHAP* (Section 2).
2. *Unified Framework*: Application of *TreeSHAP-IQ* to the broad class of any-order CII.
3. *Application*: We efficiently implement *TreeSHAP-IQ* on SOTA tree-based models, such as XGBoost, and showcase how interaction scores enrich single feature attribution measures on several benchmark datasets (Section 4).

2 Local Shapley-Based Explanations

Local Shapley-based explanations consider a model f on an n -dimensional feature space \mathcal{X} with features $N := \{1, \dots, n\}$. The goal is to explain the prediction $f(x) \in \mathbb{R}$ for a selected explanation point $x \in \mathcal{X}$ and find an additive attribution $\phi = (\phi[1], \dots, \phi[n]) \in \mathbb{R}^n$, such that $f(x) = b_0 + \sum_{i \in N} \phi[i]$, where $b_0 \in \mathbb{R}$ is the baseline prediction, i.e. the prediction of x , if no feature information is available. To compute a unique attribution score $\phi[i]$ for each feature $i \in N$, we extend the model with subsets of features $f : \mathcal{X} \times \mathcal{P}(N) \rightarrow \mathbb{R}$, where $\mathcal{P}(N)$ is the power set of N

and $f(x, T)$ refers to the prediction of f at x , if only the features in $T \subseteq N$ are known. In the following, if we omit the subset, then $T = N$, i.e. $f(x) := f(x, N)$. We further omit the explanation point x if it is clear from context, and set $f(T) := f(x, T)$ and $b_0 := f(x, \emptyset)$. The contribution of each feature $i \in N$ is then the SV (Shapley 1953)

$$\phi(f, i) := \sum_{T \subseteq N \setminus \{i\}} \frac{1}{n \cdot \binom{n-1}{|T|}} [f(T \cup \{i\}) - f(T)].$$

The SVs define the unique attribution measure satisfying the following axioms: linearity (in f), symmetry (ordering does not impact ϕ), dummy (no impact on f implies $\phi(f, i) = 0$) and efficiency $f(x) = b_0 + \sum_{i \in N} \phi(f, i)$ (Shapley 1953).

In many real-world applications, single feature importance scores are not sufficient to understand a model, where features become only meaningful when *interacting* with others. The SV does not give any information about such *interactions* between two or more features. The SII has been the first extension of the SV to interactions of feature subsets.

Definition 1 (SII, Grabisch and Roubens 1999). *The SII for an interaction $S \subseteq N$ is defined as*

$$I^{SII}(f, S) := \sum_{T \subseteq N \setminus S} \frac{1}{(n - |S| + 1) \cdot \binom{n-|S|}{|T|}} \delta_S(f, T),$$

where δ_S is the S -derivative of f for $T \subseteq N \setminus S$, i.e.

$$\delta_S(f, T) := \sum_{L \subseteq S} (-1)^{|S|-|L|} f(T \cup L).$$

The SII is the unique attribution measure that fulfills the (generalized) linearity, symmetry and dummy axiom, as well as a novel recursive axiom that links higher to lower order interactions (Grabisch and Roubens 1999). In contrast to the SV, the SII does not fulfill the (generalized) efficiency axiom, which states that the sum of interaction scores (including b_0) up to a *maximum order* s_0 equals the model prediction $f(x)$. This axiom is particularly useful in the ML context. Recently, Bordt and von Luxburg (2023) proposed a specific aggregation, known as n-SII of order s_0 , which yields a unique index that satisfies the (generalized) efficiency axiom. A more general class constitutes the CII, where it was shown that every interaction index fulfilling the linearity, symmetry and dummy axiom can be represented as a CII (Grabisch and Roubens 1999, Proposition 5). Other CIIs were proposed that introduce a unique interaction index of order s_0 and require the efficiency axiom directly, such as the STI (Sundararajan, Dhamdhere, and Agarwal 2020) or the FSI (Tsai, Yeh, and Ravikumar 2023). While the computation of the SV and SIIs are of exponential complexity, it has been shown that the complexity for the SV can be reduced to polynomial time in the case of tree-based models.

2.1 The Shapley Value for Tree Ensembles

For tree-based models the computational complexity of the SV can be drastically reduced by utilizing the additive tree structure. Furthermore, there exists a natural way to handle missing features, which can be used to define the extended model $f(x, T)$. For simplicity, we consider in the following a single decision tree, where ensembles of trees can be similarly computed due to the linearity of the SV.

$$x[i] = 0.7, x[j] = 1.6, x[k] = 0.2, f(x) = 0.65$$

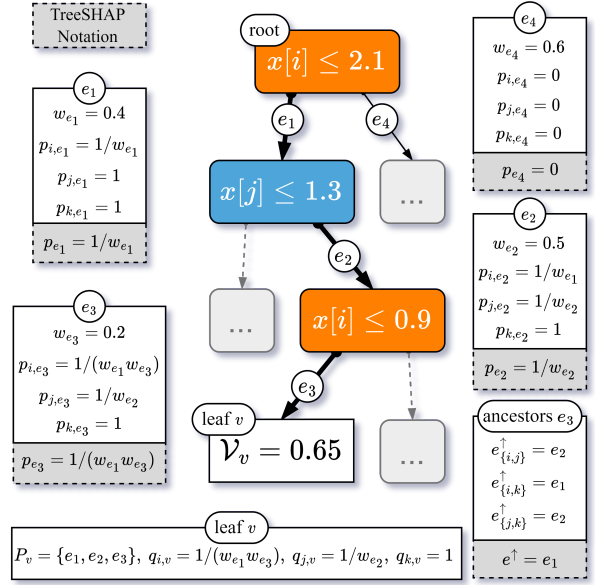


Figure 3: Notations in TreeSHAP-IQ and Linear TreeSHAP.

Notation. We consider a decision tree $\mathcal{T} = (V, E)$ as a rooted directed tree with a set of vertices V , referred to as decision nodes, and edges E . The root node is denoted as $r \in V$. Each decision node consists of a split feature $i \in N$ with a threshold value and predictions \mathcal{V}_v at the leaf nodes. For each node $v \in V$, we let P_v be the set of edges from the root node to v and $\mathcal{L}(v)$ the set of leaf nodes reachable from v , where $\mathcal{L}(\mathcal{T}) = \mathcal{L}(r)$ is the set of all leaf nodes in the tree. For every edge $e \in E$ going from u to v , we denote u as the tail of e and v as the head of e , $h(e)$. We consider a weighted tree with weights $w_e \in (0, 1)$ for every edge $e \in E$, which is defined as the proportion of observed data points at the tail of e , that split to the head of e . Additionally, we label each edge $e \in E$ with the feature associated with the tail of e , i.e. the feature that was used to split the observations on the decision node at the source of e . Further, $P_{i,v}$ and E_i are the edges in P_v and E with label $i \in N$. Our notation for decision trees is illustrated in Figure 3.

We also require polynomial arithmetic and refer to the set of polynomials with maximum degree d and coefficients in \mathbb{R} as $\mathbb{R}[x]_d$. Polynomial multiplication is denoted with \odot and division with $\lfloor \frac{a}{b} \rfloor$ or $\lfloor a/b \rfloor$. We denote with $\langle x, y \rangle$ the inner product of two vectors $x, y \in \mathbb{R}^d$ and refer to the inner product of the coefficients, if polynomials are considered.

Extended Model $f(x, T)$ for Decision Trees. A decision tree can be decomposed into distinct decision rules $R^v : \mathcal{X} \rightarrow \mathbb{R}$ for each leaf $v \in \mathcal{L}(\mathcal{T})$, which predict \mathcal{V}_v , if x reaches v and zero otherwise. Note that each R^v induces a subspace of \mathcal{X} at which the prediction of f is constant. The decision tree is thus given as $f(x) = \sum_{v \in \mathcal{L}(\mathcal{T})} R^v(x)$. We now define R_T^v , the prediction rule restricted to a set of *active* features $T \subseteq N$, where the remaining are con-

sidered to be unknown. If the split feature is unknown, we split based on the weights w_e , which is a common practice (Yu et al. 2022). If all features are unknown, we define $R_\emptyset^v := R_\emptyset^v(x) := \mathcal{V}_v \prod_{e \in P_v} w_e$. When adding feature $i \in N$ to the active set, the product of associated weights is replaced by the split criterion. This is formalized as a recursive property $R_{T \cup \{i\}}^v = q_{i,v}(x) R_T^v(x)$, where $q_{i,v}$ is the marginal effect of adding i to the active set. To define $q_{i,v}$, we let $x \in \pi_i(R^v)$, if $x[i]$ satisfies each split criterion regarding features i in the path of v , i.e. $\pi_i(R^v)$ is the region of feature i in the induced subspace of \mathcal{X} by R^v . For v the marginal effect of adding feature $i \in N$ is then defined as

$$q_{i,v}(x) := \mathbf{1}(x \in \pi_i(R^v)) \prod_{e \in P_{i,v}} \frac{1}{w_e}, \quad (1)$$

where $\mathbf{1}(\cdot)$ is the indicator function. Furthermore, for $P_{i,v} = \emptyset$ we define $q_{i,v} = 1$. The restricted rule is thus defined as

$$R_T^v(x) := R_\emptyset^v \prod_{j \in T} q_{j,v}(x). \quad (2)$$

For a tree \mathcal{T}_f and $T \subseteq N$ the restricted model at x is then

$$f(T) := f(x, T) := \sum_{v \in \mathcal{L}(\mathcal{T}_f)} R_T^v(x).$$

In the following, we omit the argument x in the notation. We proceed to compute the SV of $f(T)$, which is known as *path dependent TreeSHAP* (Lundberg et al. 2020).

Linear TreeSHAP. TreeSHAP exploits the tree structure to compute the SV in polynomial time (Lundberg et al. 2020). Linear TreeSHAP improved this computation and provided a theoretical framework by linking the computation to polynomial arithmetic (Yu et al. 2022). Plugging (2) into the definition of the SV and using the fact that $q_{j,v} - 1 = 0$, if a feature does not appear in the path, yields

$$\phi(R^v, i) = (q_{i,v} - 1) \sum_{T \subseteq \mathcal{F}(R^v) \setminus \{i\}} \frac{R_\emptyset^v}{n \cdot \binom{n-1}{|T|}} \prod_{j \in T} q_{j,v}, \quad (3)$$

where $\mathcal{F}(R^v)$ is the set of all features that appear in R^v . It was shown that this sum can be efficiently stored using the coefficients of a specific polynomial.

Definition 2 (Summary Polynomial (SP), Yu et al. 2022). *The SP of leaf node v is $G_v^{SP}(y) := R_\emptyset^v \prod_{j \in \mathcal{F}(R^v)} (q_{j,v} + y)$.*

For feature $i \in N$, $R_\emptyset^v \sum_{S \subseteq \mathcal{F}(R) \setminus \{i\}} \prod_{j \in S} q_{j,v}$ is the coefficient of y^{d-k-1} in $\frac{G}{q_{i,v} + y}$, where $d := |\mathcal{F}(R^v)|$ is the number of features in each path. Note that this corresponds to the non-weighted terms in the sum of (3) for $k = 0, \dots, d-1$. The SV of a single decision rule can thus be represented as

$$\phi(R^v, i) = (q_{i,v} - 1) \psi \left(\left[\frac{G_v}{q_{i,v} + y} \right] \right), \quad (4)$$

where $\psi : \mathbb{R}[x]_d \rightarrow \mathbb{R}$ is a function that properly weights the coefficients, such that it corresponds to the sum in (3). It

is formally defined (Yu et al. 2022) as

$$\psi_d(A) := \frac{\langle A, B_d \rangle}{d+1} \text{ with } B_d(y) := \sum_{k=0}^d \binom{d}{k}^{-1} y^k. \quad (5)$$

We write $\psi(A) = \psi_d(A)$, where d is the degree of A . It was then shown that ψ is additive and scale invariant.

Proposition 1 (Yu et al. 2022). *For ψ and $p, q \in \mathbb{R}[x]_d$,*

$$\psi_d(p+q) = \psi_d(p) + \psi_d(q) \text{ and } \psi(p \odot (1+y)^k) = \psi(p).$$

Using (4) based on leaf nodes, a representation of the SV in terms of edges is presented, which is explicitly computed by recursively traversing the tree. For this representation, the SP is extended to every edge in the path of $P_{i,v}$ as

$$G_u := \bigoplus_{v \in \mathcal{L}(u)} G_v \text{ with } G^1 \oplus G^2 := G^1 + G^2 \odot (1+y)^{d_1-d_2},$$

where the order is such that $d_1 > d_2$, i.e. \oplus is an operation on the set of polynomials $\oplus : \mathbb{R}[x]_{d_1} \times \mathbb{R}[x]_{d_2} \rightarrow \mathbb{R}[x]_{\max(d_1, d_2)}$ that sums the polynomial while scaling them to the same degree. Note that due to the properties of ψ , we have $\psi(G_u) = \sum_{v \in \mathcal{L}(u)} \psi(G_v)$. For edge $e \in E$ and its feature i , we further introduce the inter-path value of $q_{i,v}$ as

$$p_e := \mathbf{1}(x \in \pi_{h(e)}) \prod_{e' \in P_{i, h(e)}} \frac{1}{w_{e'}}.$$

Note that $p_{e^*} = q_{i,v}$ if e^* is the last edge in $P_{i,v}$. An edge-based representation of the SV is then provided.

Theorem 1 (Yu et al. 2022). *Let $i \in N$ and denote for e the closest ancestor in the set E_i by e^\uparrow , where $e^\uparrow = \perp$ and $p_{i, \perp} = 1$ in case it does not exist. Then,*

$$\begin{aligned} \phi(f, i) &= \sum_{e \in E_i} (p_e - 1) \psi \left(\left[\frac{G_{h(e)}}{y + p_e} \right] \right) \\ &\quad - (p_{e^\uparrow} - 1) \psi \left(\left[\frac{G_{h(e)} \odot (y+1)^{d_{e^\uparrow} - d_e}}{y + p_{e^\uparrow}} \right] \right). \end{aligned}$$

Using this edge-based representation, Linear TreeSHAP computes the SV by traversing once through the tree. To improve efficiency, the SP is stored in a multipoint interpolation form. For more details, we refer to Appendix B.

3 TreeSHAP-IQ: Computation of Local Shapley Interactions for Tree Ensembles

Computing the exact SV for tree ensembles can reliably quantify the impact of single features on the model's predictions. However, in many applications, certain features become only meaningful when *interacting* with other features. In this case, the SV is not sufficient to understand how the model predicts, and more complex explanations in terms of Shapley interactions are necessary. In the following, we propose TreeSHAP Interaction Quantification (TreeSHAP-IQ), an efficient algorithm for computing any-order SII scores, which follows naturally by extending the SP to interactions.

TreeSHAP-IQ can further be applied to the broad class of CII (Grabisch and Roubens 1999), which we briefly discuss in Section 3.2. All proofs are deferred to Appendix A.

3.1 Theoretical Foundation of TreeSHAP-IQ

We now present the theoretical foundation of TreeSHAP-IQ. The notations in this section extend on Linear TreeSHAP (Yu et al. 2022) and are illustrated in Figure 3. We compute the S-derivative for R^v and $T \subseteq N \setminus S$ as

$$\delta_S(R^v, T) = R_T^v \sum_{L \subset S} (-1)^{|S|-|L|} \prod_{j \in L} q_{j,v}, \quad (6)$$

which follows from (2) and the recursive property. We thus represent the SII for a single decision rule as follows.

Proposition 2. For a leaf v in \mathcal{T}_f , it holds $I^{SII}(R^v, S) =$

$$\left(\sum_{L \subset S} (-1)^{|S|-|L|} \prod_{j \in L} q_{j,v} \right) \psi \left(\left\lfloor \frac{G_v}{\prod_{j \in S} (q_{j,v} + y)} \right\rfloor \right).$$

Proposition 2 yields a compact representation in terms of leaf nodes and decision rules, which reduces to the representation of (4) for single feature subsets. Similar to Linear TreeSHAP, the representation of SII in terms of leaf nodes is not suitable for efficient computation. We thus again establish an edge-based representation, similar to Theorem 1. By Proposition 2, the computation of an interaction for a subset $S \subset N$ requires knowledge of all $q_{i,v}$ with $i \in S$, which have to be tracked during the traversal of the tree. We thus first extend the inter-path values p_e to every feature as

$$p_{i,e} := \mathbf{1}(x \in \pi_{i,h(e)}) \prod_{e' \in P_{i,h(e)}} \frac{1}{w_{e'}},$$

where $x \in \pi_{i,u}$ if $x[i]$ satisfies each decision criterion in $P_{i,u}$. Note that $p_{j,e} = p_e$ and $\pi_{j,h(e)} = \pi_{h(e)}$, if j is the label of e . Our goal in the following is to provide an algorithm similar to Linear TreeSHAP that traverses the decision tree once and recursively computes the interaction scores. The SP thereby remains unchanged, but we introduce further polynomials of order $|S|$ to efficiently maintain the sum as well as the denominator in Proposition 2.

Definition 3 (Interaction Polynomial (IP)). The IP of $S \subset N$ and edge e is $H_{S,e}^{IP}(y) := \prod_{j \in S} (p_{j,e} - y)$.

Note that the coefficient of y^k in $H_{S,e}$ is exactly $\sum_{L \subset S} (-1)^{|S|-|L|} \prod_{j \in L} p_{j,e}$ for $k = 0, \dots, |S|$. Therefore, the sum of the coefficients of the IP equals the sum in (6). We thus define the coefficient sum.

Definition 4 (Coefficient sum κ). We define the function $\kappa_d : \mathbb{R}[x]_d \rightarrow \mathbb{R}$ as $\kappa_d(A) := \langle A, y^d + \dots + y + 1 \rangle$. We write $\kappa(p) = \kappa_d(p)$, where d is the degree of p .

Applying κ to the IP yields the following properties.

Proposition 3. For the sum of coefficients of the IP, it holds

$$\kappa(H_{S,e}^{IP}) = \sum_{L \subset S} (-1)^{|S|-|L|} \prod_{j \in L} p_{j,e}. \quad (7)$$

If there exists $j \in S$ with $p_{j,e} = 1$, then $\kappa(H_{S,e}^{IP}) = 0$.

Proposition 3 shows that $\kappa(H_{S,e}^{IP})$ corresponds to the edge-based representation of the sum in Proposition 2. If e is

the last edge in $P_{j,v}$, then $p_{j,e} = q_{j,v}$ for all $j \in N$ and thus $\kappa(H_{S,e}^{IP})$ retrieves the sum in Proposition 2. Furthermore, if $p_{j,e} = 1$, then it is intuitive that all inter-path contributions with $j \in S$ are zero, since j does not impact the model's prediction in this part of the tree. This property allows us to update interaction scores only if all features of the subset have occurred in the path. We further describe the quotient in Proposition 2 using another polynomial of order $|S|$.

Definition 5 (Quotient Polynomial (QP)). The QP of $S \subset N$ and edge e is $H_{S,e}^{QP}(y) := \prod_{j \in S} (p_{j,e} + y)$.

If e_S^* is the last edge in P_v of leaf node v that contains any feature of S , then $p_{j,e_S^*} = q_{j,v}$ for every $j \in S$ and hence we can rewrite Proposition 2 using Proposition 3 as

$$I^{SII}(R^v, S) = \kappa(H_{S,e_S^*}^{IP}) \psi \left(\left\lfloor G_v / H_{S,e_S^*}^{QP} \right\rfloor \right). \quad (8)$$

Clearly, Proposition 2 reduces to (4) for the case of the SV. In contrast to the SV, the edge-based computation includes all inter-path values of $p_{j,e}$ with $j \in S$. To extend Theorem 2, we therefore need to extend the notion of ancestor edges to ancestors with respect to a subset $S \subset N$.

Proposition 4. For a decision rule R^v of a leaf node v and a subset $S \subset N$, let $P_{S,v} := \bigcup_{i \in S} P_{i,v}$ and e_S^\uparrow as the closest ancestor of e in $P_{S,v}$. The SII of R^v is then given by

$$I^{SII}(R^v, S) = \sum_{e \in P_{S,v}} \kappa(H_{S,e}^{IP}) \psi \left(\left\lfloor \frac{G_v \odot (y+1)^{d_e - d_v}}{H_{S,e}^{QP}} \right\rfloor \right) - \kappa(H_{S,e_S^\uparrow}^{IP}) \psi \left(\left\lfloor \frac{G_v \odot (y+1)^{d_{e_S^\uparrow} - d_v}}{H_{S,e_S^\uparrow}^{QP}} \right\rfloor \right).$$

Using Proposition 4, we can state our main theorem.

Theorem 2. For $S \subset N$, let $E_S := \bigcup_{i \in S} E_i$ be the set of edges that split on any feature in S , and denote the closest ancestor of e in $P_{S,v}$ as e_S^\uparrow . The SII is then computed as

$$I^{SII}(f, S) = \sum_{e \in E_S} \kappa(H_{S,e}^{IP}) \psi \left(\left\lfloor \frac{G_{h(e)}}{H_{S,e}^{QP}} \right\rfloor \right) - \kappa(H_{S,e_S^\uparrow}^{IP}) \psi \left(\left\lfloor \frac{G_{h(e)} \odot (y+1)^{d_{e_S^\uparrow} - d_e}}{H_{S,e_S^\uparrow}^{QP}} \right\rfloor \right).$$

Note that for $S = \{i\}$, Theorem 2 reduces to Theorem 1.

Implementation of TreeSHAP-IQ. Theorem 2 allows for an efficient computation of the SII, with the SP being handled alike to Linear TreeSHAP. The IQ and QP are updated for each interaction subset that contains the feature of e . We again use the multipoint interpolation form to store and update the polynomials G_v , $H_{S,e}^{IP}$, and $H_{S,e}^{QP}$. TreeSHAP-IQ traverses the decision tree once for every explanation point. At each edge (decision node), TreeSHAP-IQ updates all interactions that contain the currently encountered feature, $\binom{n-1}{|S|-1}$ in total. However, the update can be restricted to those interactions, where all features have been observed in the path. We refer to Appendix B for more details.

Complexity of TreeSHAP-IQ Consider m explanation points, $\ell_{\mathcal{T}} := |\mathcal{L}(\mathcal{T})|$ as the number of leaves and d_{\max} as the maximum depth of the tree.

Linear TreeSHAP has a computational complexity of $\mathcal{O}(m \cdot \ell_{\mathcal{T}} \cdot d_{\max})$ and storage complexity of $\mathcal{O}(d_{\max}^2 + n)$ (Yu et al. 2022). We now consider the complexity of TreeSHAP-IQ, if all interactions of order $s := |S|$ are computed. In contrast to Linear TreeSHAP and the SP, where only the current feature value has to be updated, TreeSHAP-IQ needs to update the IP, the QP and the interaction scores for all interaction subsets that contain the currently observed feature. This increases the computational complexity by a factor of $\binom{n-1}{s-1}$. Furthermore, all interaction scores have to be stored, requiring storage of $\binom{n}{s}$. To store the IQ and QP, we require further a storage capacity of $\mathcal{O}(d_{\max}^2 \cdot \binom{n}{s})$. The computational complexity is thus summarized as follows.

TreeSHAP-IQ complexity for the SII of order s

Computational Complexity	Storage Complexity
$\mathcal{O}\left(m \cdot \ell_{\mathcal{T}} \cdot d_{\max} \cdot \binom{n-1}{s-1}\right)$	$\mathcal{O}\left(d_{\max}^2 \cdot \binom{n}{s}\right)$

For the computation of the SV, the computational complexity of TreeSHAP-IQ is similar to Linear TreeSHAP. The storage capacity is increased by a factor n , as we store the IP and QP for every feature. Moreover, for pairwise interactions, TreeSHAP-IQ mirrors the complexity of the computation proposed by Lundberg et al. (2020) using Linear TreeSHAP. However, our method distinguishes itself by relying on a single initialization of the tree parameters.

3.2 Extending TreeSHAP-IQ to General CII

TreeSHAP-IQ can be extended to the broad class of CII. A CII is defined as $I^{\text{CII}}(f, S) := \sum_{T \subseteq N \setminus S} w_s^{\text{CII}}(|T|) \delta_S(f, T)$ with non-negative weights w_s^{CII} that depend on the interaction order $s := |S|$ (Grabisch and Roubens 1999; Fumagalli et al. 2023). This includes other approaches of extending the SV to interactions, such as STI and FSI, as well as Banzhaf interactions (Patel, Strobel, and Zick 2021). Observe from the proofs, that different weights in CII solely impact the SP, and in particular ψ . To extend the SP for CII, we let $d := |\mathcal{F}(R^v)|$ and scale G_v to the degree of n , which does not impact ψ due to the scale invariance. We then observe

$$\psi \left(\left[\frac{G_v \odot (1+y)^{n-d}}{\prod_{j \in S} (1+q_{j,v})} \right] \right) = R_{\emptyset}^v \sum_{T \subseteq N \setminus S} w_s^{\text{SII}}(t) \prod_{j \in S} q_{j,v},$$

where $w_s^{\text{SII}}(t)$ is the CII weight for SII, cf. Definition 1. Recall from (5) that these weights are retrieved from the polynomial B_{n-s} . Thus, we generalize $\psi_d^{\text{CII}}: \mathbb{R}[x]_d \rightarrow \mathbb{R}$ to

$$\psi_d^{\text{CII}}(A) := \langle A, W_d \rangle \text{ with } W_d^{\text{CII}}(y) := \sum_{k=0}^d w_{n-d}^{\text{CII}}(k) y^k.$$

If G_v is scaled to degree n , then ψ_d^{CII} is always evaluated with a polynomial of degree $d = n - |S|$. Further, note that for SII, we have $\psi_d(A) \equiv \psi_d^{\text{SII}}(A)$, where the quotient $d + 1$ is included in the weights, i.e. $W_d^{\text{SII}}(y) = B_d / (d + 1)$.

Datasets	# Instances	# Features	Target	Speed-Up
<i>Credit</i>	1 000	20	{0, 1}	$\sim 10^4$
<i>Bank</i>	45 211	16	{0, 1}	$\sim 10^3$
<i>Adult</i>	45 222	14	{0, 1}	$\sim 10^3$
<i>Bike</i>	17 379	12	\mathbb{R}	$\sim 10^1$
<i>COMPAS</i>	6 172	11	{0, 1}	$\sim 10^2$
<i>Titanic</i>	891	9	{0, 1}	$\sim 10^1$
<i>California</i>	20 640	8	\mathbb{R}	~ 1

Table 1: Overview of datasets and speed-up compared to a naive computation

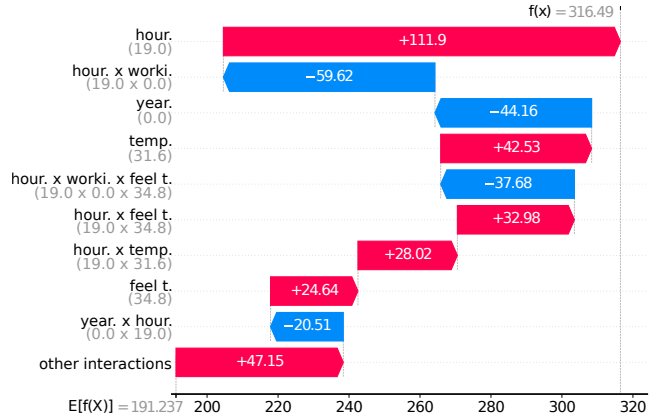


Figure 4: Waterfall chart for n-SII scores with $s_0 = 3$ and a prediction of the *Bike* regression dataset.

Implementation of CII in TreeSHAP-IQ Using ψ^{CII} , any CII can be computed by TreeSHAP-IQ. In contrast to $\psi \equiv \psi^{\text{SII}}$, the scale invariance does not hold for CII. Therefore, the SP cannot be reduced to the degree $d := |\mathcal{F}(R^v)|$. However, if we maintain the SP at the maximum degree n , then all previous results apply. If the SP is stored in multipoint interpolation form, then this merely requires a multiplication with the corresponding term of $(y + 1)^{n-d}$, which can be efficiently precalculated. Thus, the computational complexity is not affected by this extension. Provided $d_{\max} \geq n$, the storage complexity is not affected either.

4 Experiments

We apply TreeSHAP-IQ¹ on XGBoost (XBG) (Chen and Guestrin 2016), gradient-boosted trees (GBTs), random forest (RF), and decision tree (DT) algorithms on the *German Credit* (Hofmann 1994), *Bank* (Moro, Cortez, and Laureano 2011), *Adult Census* (Kohavi 1996), *Bike* (Fanace-T and Gama 2014), *COMPAS* (Angwin et al. 2016), *Titanic* (Dawson 1995), and *California* (Kelley Pace and Barry 1997) datasets, see Table 1. For further experimental results, including a run-time analysis and detailed information on the datasets, models, and pre-processing steps, we refer to Appendix C. We compute additive interactions for single predictions using TreeSHAP-IQ with n-SII of different order.

¹All experimental code and the technical appendix can be found at: github.com/mmschlk/TreeSHAP-IQ.

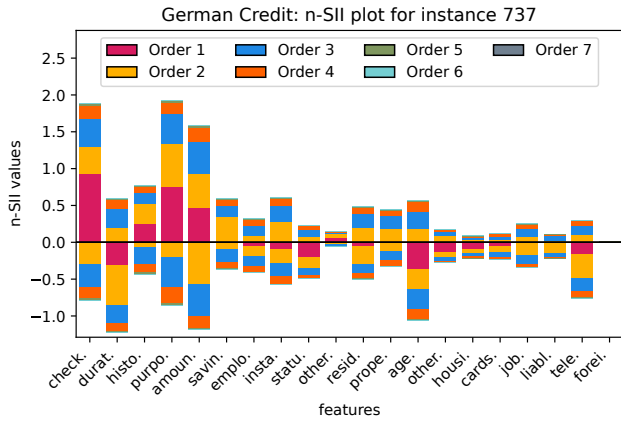


Figure 5: Visualization of positive and negative n-SII scores per feature with $s_0 = 7$ for an observation in *German Credit*.

TreeSHAP-IQ Reveals Intricate Feature Interactions.

Using TreeSHAP-IQ, we examine the model’s prediction based on higher order interaction effects. We distinguish n-SII scores that positively (red) and negatively (blue) impact the prediction. In Figure 1, we visualize n-SII with $s_0 = 2$. The width of the network vertices (order 1) and the network edges (order 2) describes the absolute value of the corresponding n-SII scores. We observe that there exist features that strongly impact the prediction individually, such as the information about a non-existing *checking account*. However, the present *credit amount* strongly impacts the prediction only in interaction with the given *installment rate* (positively) and *duration* (negatively).

The force plots in Figure 2 illustrate how the additive local explanations change, if higher order interactions are considered. We consider the n-SII scores for $s_0 = 1, 2, 3$ for the *California* housing dataset and an XGBoost regressor. The force plot displays the positive and negative interaction scores starting from the predicted value to the left and right, respectively, sorted by their absolute value. We observe that individual feature effects, such as *Longitude*, reduce when higher order interactions are considered. The interaction of *Longitude* and *Latitude* reveals the importance of the geographic location of this instance.

The waterfall chart in Figure 4 displays the explanations of n-SII with order $s_0 = 3$ for an instance in the *bike* dataset. For this instance, it can be seen that the interaction of the evening *hour* with a non-*working* day affects the prediction negatively, whereas the interaction with both *temperature* features contribute positively.

n-SII Plots Quantify Interactions of Each Feature.

To assess the strength of interaction per individual feature, we utilize the visualization of n-SII values presented by Bordt and von Luxburg (2023). We compute exact n-SII scores up to order $s_0 = 7$ for the *German Credit* dataset. The positive and negative interactions are distributed equally onto each feature in the subset and displayed on the positive and negative axes, respectively. The sum of all stacked bars results in the SV of each feature (Bordt and von Luxburg 2023). In

Figure 5, we observe that the interaction effects diminish at order 5, with interactions of orders 6 and 7 being virtually absent. Assuming that interactions decay with higher order, this visualization can be used to find the maximum order to explain the prediction (i.e. $s_0 = 5$ from Figure 5).

5 Limitations

TreeSHAP-IQ applies to the broad class of CII, provided that its representation in terms of a weighted sum of discrete derivatives is known. For FSI, this representation is only explicitly known for top-order interactions (Tsai, Yeh, and Ravikumar 2023), as FSI is motivated as a solution to a constrained weighted least square problem. Similar to the SV, Shapley interactions strongly rely on how absent features are modeled. In our work, we considered the *path dependent* feature perturbation (Lundberg et al. 2020), which is linked to the *observational* approach (Chen et al. 2020). The *interventional* approach (Lundberg et al. 2020) can be computed with TreeSHAP-IQ, akin to TreeSHAP, but similarly increases the computational complexity by the number of samples used in the background dataset. In this case, more efficient variants should be used instead (Zern, Broelemann, and Kasneci 2023). Both paradigms yield different explanations, where the appropriate choice should be carefully done depending on the application (Chen et al. 2020).

6 Conclusion and Future Work

We presented TreeSHAP-IQ, an efficient method to compute any-order additive Shapley interactions that locally explain single predictions for general ensembles of trees. Akin to SOTA Linear TreeSHAP (Yu et al. 2022), our algorithm is based on a solid theoretical foundation that exploits polynomial arithmetic. We applied TreeSHAP-IQ on SOTA ML models, such as XGBoost (Chen and Guestrin 2016), and several benchmark datasets. We demonstrated that TreeSHAP-IQ reveals intricate feature interactions, which enrich Shapley-based feature attribution.

Utilizing well-known visualization and aggregation techniques from machine learning (Lundberg and Lee 2017; Bordt and von Luxburg 2023) and statistics (Inglis, Parnell, and Hurley 2022) we presented these scores in a manner that is easily understandable and interpretable. While interactions are widely studied in statistics, explaining local predictions using interaction scores, in particular with Shapley-based interactions, is an emerging line of research in the field of XAI. Due to the exponentially increasing number of interactions, we provided intuitive visualizations to present TreeSHAP-IQ scores to practitioners. Nevertheless, it would be beneficial to explore further human-centered post-processing techniques and visualizations, as well as rigorously evaluate the explanatory capabilities of TreeSHAP-IQ with user studies, especially to validate quantitatively that the user’s understanding increases when higher order explanations are presented. Additionally, the n-SII scores define a local generalized additive model (GAM) (Bordt and von Luxburg 2023) that could be further linked to functional decomposition (Hiabu, Meyer, and Wright 2023).

Acknowledgements

We sincerely thank the anonymous reviewers for their work and helpful comments. We gratefully acknowledge funding by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation): TRR 318/1 2021 – 438445824.

References

- Adadi, A.; and Berrada, M. 2018. Peeking Inside the Black-Box: A Survey on Explainable Artificial Intelligence (XAI). *IEEE Access*, 6: 52138–52160.
- Angwin, J.; Larson, J.; Mattu, S.; and Kirchner, L. 2016. Machine Bias. There’s software used across the country to predict future criminals. And it’s biased against blacks.
- Bordt, S.; and von Luxburg, U. 2023. From Shapley Values to Generalized Additive Models and back. In *International Conference on Artificial Intelligence and Statistics (AISTATS 2023)*, volume 206 of *Proceedings of Machine Learning Research*, 709–745. PMLR.
- Boulesteix, A.-L.; Bender, A.; Lorenzo Bermejo, J.; and Strobl, C. 2012. Random forest Gini importance favours SNPs with large minor allele frequency: impact, sources and recommendations. *Briefings in Bioinformatics*, 13(3): 292–304.
- Casalicchio, G.; Molnar, C.; and Bischl, B. 2018. Visualizing the Feature Importance for Black Box Models. In *Machine Learning and Knowledge Discovery in Databases - European Conference (ECML PKDD 2018)*, volume 11051 of *Lecture Notes in Computer Science*, 655–670. Springer.
- Castro, J.; Gómez, D.; Molina, E.; and Tejada, J. 2017. Improving polynomial estimation of the Shapley value by stratified random sampling with optimum allocation. *Computers & Operations Research*, 82: 180–188.
- Castro, J.; Gómez, D.; and Tejada, J. 2009. Polynomial calculation of the Shapley value based on sampling. *Computers & Operations Research*, 36(5): 1726–1730.
- Chen, H.; Covert, I.; Lundberg, S.; et al. 2023. Algorithms to estimate Shapley value feature attributions. *Nature Machine Intelligence*, 5: 590–601.
- Chen, H.; Janizek, J. D.; Lundberg, S. M.; and Lee, S. 2020. True to the Model or True to the Data? arXiv:2006.16234.
- Chen, T.; and Guestrin, C. 2016. XGBoost: A Scalable Tree Boosting System. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (SIGKDD 2016)*, 785–794. ACM.
- Covert, I.; and Lee, S. 2021. Improving KernelSHAP: Practical Shapley Value Estimation Using Linear Regression. In *The 24th International Conference on Artificial Intelligence and Statistics (AISTATS 2021)*, volume 130 of *Proceedings of Machine Learning Research*, 3457–3465. PMLR.
- Covert, I.; Lundberg, S. M.; and Lee, S. 2020. Understanding Global Feature Contributions With Additive Importance Measures. In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020 (NeurIPS 2020)*.
- Dawson, R. J. M. 1995. The “Unusual Episode” Data Revisited. *Journal of Statistics Education*, 3(3).
- Fanaee-T, H.; and Gama, J. 2014. Event Labeling Combining Ensemble Detectors and Background Knowledge. *Progress in Artificial Intelligence*, 2(2): 113–127.
- Feurer, M.; van Rijn, J. N.; Kadra, A.; Gijsbers, P.; Mallik, N.; Ravi, S.; Mueller, A.; Vanschoren, J.; and Hutter, F. 2020. OpenML-Python: an extensible Python API for OpenML. arXiv:1911.02490.
- Friedman, J. H. 2001. Greedy function approximation: A gradient boosting machine. *The Annals of Statistics*, 29(5): 1189–1232.
- Fumagalli, F.; Muschalik, M.; Kolpaczki, P.; Hüllermeier, E.; and Hammer, B. 2023. SHAP-IQ: Unified Approximation of any-order Shapley Interactions. arXiv:2303.01179.
- Grabisch, M.; and Roubens, M. 1999. An axiomatic approach to the concept of interaction among players in cooperative games. *International Journal of Game Theory*, 28(4): 547–565.
- Harris, C.; Pymar, R.; and Rowat, C. 2022. Joint Shapley values: a measure of joint feature importance. In *The Tenth International Conference on Learning Representations (ICLR 2022)*. OpenReview.net.
- Herbinger, J.; Bischl, B.; and Casalicchio, G. 2023. Decomposing Global Feature Effects Based on Feature Interactions. arXiv:2306.00541.
- Hiabu, M.; Meyer, J. T.; and Wright, M. N. 2023. Unifying local and global model explanations by functional decomposition of low dimensional structures. In *International Conference on Artificial Intelligence and Statistics (AISTATS 2023)*, volume 206 of *Proceedings of Machine Learning Research*, 7040–7060. PMLR.
- Hofmann, H. 1994. Statlog (German Credit Data). UCI Machine Learning Repository. DOI: 10.24432/C5NC77.
- Hooker, G. 2004. Discovering additive structure in black box functions. In Kim, W.; Kohavi, R.; Gehrke, J.; and DuMouchel, W., eds., *Proceedings of the Tenth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (SIGKDD 2004)*, 575–580. ACM.
- Hooker, G. 2007. Generalized Functional ANOVA Diagnostics for High-Dimensional Functions of Dependent Variables. *Journal of Computational and Graphical Statistics*, 16(3): 709–732.
- Inglis, A.; Parnell, A.; and Hurley, C. B. 2022. Visualizing Variable Importance and Variable Interaction Effects in Machine Learning Models. *Journal of Computational and Graphical Statistics*, 31(3): 766–778.
- Jethani, N.; Sudarshan, M.; Covert, I. C.; Lee, S.; and Ranganath, R. 2022. FastSHAP: Real-Time Shapley Value Estimation. In *The Tenth International Conference on Learning Representations (ICLR 2022)*. OpenReview.net.
- Ke, G.; Meng, Q.; Finley, T.; Wang, T.; Chen, W.; Ma, W.; Ye, Q.; and Liu, T. 2017. LightGBM: A Highly Efficient Gradient Boosting Decision Tree. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017 (NeurIPS 2017)*, 3146–3154.

- Kelley Pace, R.; and Barry, R. 1997. Sparse spatial autoregressions. *Statistics & Probability Letters*, 33(3): 291–297.
- Kohavi, R. 1996. Scaling up the Accuracy of Naive-Bayes Classifiers: A Decision-Tree Hybrid. In *Proceedings of International Conference on Knowledge Discovery and Data Mining (KDD 1996)*, 202–207.
- Kolpaczki, P.; Bengs, V.; Muschalik, M.; and Hüllermeier, E. 2023. Approximating the Shapley Value without Marginal Contributions. arXiv:2302.00736.
- Kumar, I.; Scheidegger, C.; Venkatasubramanian, S.; and Friedler, S. A. 2021. Shapley Residuals: Quantifying the limits of the Shapley value for explanations. In *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021 NeurIPS 2021*, 26598–26608.
- Kumar, I. E.; Venkatasubramanian, S.; Scheidegger, C.; and Friedler, S. A. 2020. Problems with Shapley-value-based explanations as feature importance measures. In *Proceedings of the 37th International Conference on Machine Learning (ICML 2020)*, volume 119 of *Proceedings of Machine Learning Research*, 5491–5500. PMLR.
- Lengerich, B. J.; Tan, S.; Chang, C.; Hooker, G.; and Caruana, R. 2020. Purifying Interaction Effects with the Functional ANOVA: An Efficient Algorithm for Recovering Identifiable Additive Models. In *The 23rd International Conference on Artificial Intelligence and Statistics (AISTATS 2020)*, volume 108 of *Proceedings of Machine Learning Research*, 2402–2412. PMLR.
- Lundberg, S. M.; Erion, G. G.; Chen, H.; DeGrave, A. J.; Prutkin, J. M.; Nair, B.; Katz, R.; Himmelfarb, J.; Bansal, N.; and Lee, S. 2020. From local explanations to global understanding with explainable AI for trees. *Nature Machine Intelligence*, 2(1): 56–67.
- Lundberg, S. M.; and Lee, S. 2017. A Unified Approach to Interpreting Model Predictions. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, NeurIPS 2017*, 4765–4774.
- Lunetta, K. L.; Hayward, L. B.; Segal, J.; and Van Eerdewegh, P. 2004. Screening large-scale association study data: exploiting interactions using random forests. *BMC Genetics*, 5: 32.
- Moro, S.; Cortez, P.; and Laureano, R. 2011. Using Data Mining for Bank Direct Marketing: An Application of the CRISP-DM Methodology. In *Proceedings of the European Simulation and Modelling Conference (ESM 2011)*.
- Patel, N.; Strobel, M.; and Zick, Y. 2021. High Dimensional Model Explanations: An Axiomatic Approach. In *2021 ACM Conference on Fairness, Accountability, and Transparency, Virtual Event (FAccT 2021)*, 401–411. ACM.
- Pedregosa, F.; Varoquaux, G.; Gramfort, A.; Michel, V.; Thirion, B.; Grisel, O.; Blondel, M.; Prettenhofer, P.; Weiss, R.; Dubourg, V.; VanderPlas, J.; Passos, A.; Cournapeau, D.; Brucher, M.; Perrot, M.; and Duchesnay, E. 2011. Scikit-learn: Machine Learning in Python. *Journal of Machine Learning Research*, 12: 2825–2830.
- Shapley, L. S. 1953. A Value for n-Person Games. In *Contributions to the Theory of Games (AM-28), Volume II*, 307–318. Princeton University Press.
- Shwartz-Ziv, R.; and Armon, A. 2022. Tabular data: Deep learning is not all you need. *Information Fusion*, 81: 84–90.
- Slack, D.; Hilgard, S.; Jia, E.; Singh, S.; and Lakkaraju, H. 2020. Fooling LIME and SHAP: Adversarial Attacks on Post hoc Explanation Methods. In *AAAI/ACM Conference on AI, Ethics, and Society (AIES 2020)*, 180–186. ACM.
- Strumbelj, E.; and Kononenko, I. 2014. Explaining prediction models and individual predictions with feature contributions. *Knowledge and Information Systems*, 41(3): 647–665.
- Sundararajan, M.; Dhamdhere, K.; and Agarwal, A. 2020. The Shapley Taylor Interaction Index. In *Proceedings of the 37th International Conference on Machine Learning (ICML 2020)*, volume 119 of *Proceedings of Machine Learning Research*, 9259–9268. PMLR.
- Sundararajan, M.; and Najmi, A. 2020. The Many Shapley Values for Model Explanation. In *Proceedings of the 37th International Conference on Machine Learning (ICML 2020)*, volume 119 of *Proceedings of Machine Learning Research*, 9269–9278. PMLR.
- Tsai, C.; Yeh, C.; and Ravikumar, P. 2023. Faith-Shap: The Faithful Shapley Interaction Index. *Journal of Machine Learning Research*, 24(94): 1–42.
- Tsang, M.; Rambhatla, S.; and Liu, Y. 2020. How does This Interaction Affect Me? Interpretable Attribution for Feature Interactions. In *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems (NeurIPS 2020)*, 6147–6159.
- Winham, S. J.; Colby, C. L.; Freimuth, R. R.; Wang, X.; de Andrade, M.; Huebner, M.; and Biernacka, J. M. 2012. SNP interaction detection with Random Forests in high-dimensional genetic data. *BMC Bioinformatics*, 13: 164.
- Wright, M. N.; Ziegler, A.; and König, I. R. 2016. Do little interactions get lost in dark random forests? *BMC Bioinformatics*, 17: 145.
- Yang, J. 2021. Fast TreeSHAP: Accelerating SHAP Value Computation for Trees. arXiv:2109.09847.
- Yu, P.; Bifet, A.; Read, J.; and Xu, C. 2022. Linear tree shap. In *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, (NeurIPS 2022)*.
- Zern, A.; Broelemann, K.; and Kasneci, G. 2023. Interventional SHAP Values and Interaction Values for Piecewise Linear Regression Trees. In *Thirty-Seventh AAAI Conference on Artificial Intelligence, (AAAI 2023)*, 11164–11173. AAAI Press.
- Zhang, H.; Xie, Y.; Zheng, L.; Zhang, D.; and Zhang, Q. 2021. Interpreting Multivariate Shapley Interactions in DNNs. In *Thirty-Fifth AAAI Conference on Artificial Intelligence, (AAAI 2021)*, 10877–10886. AAAI Press.

A Proofs

A.1 Proof of Proposition 2

For a leaf v in \mathcal{T}_f , it holds $I^{\text{SII}}(R^v, S) =$

$$\left(\sum_{L \subseteq S} (-1)^{|S|-|L|} \prod_{j \in L} q_{j,v} \right) \psi \left(\left\lfloor \frac{G_v}{\prod_{j \in S} (q_{j,v} + y)} \right\rfloor \right).$$

Proof. Let R^v be a decision rule of leaf node $v \in \mathcal{L}(\mathcal{T})$. For the S-derivative it follows by the recursive property and (3) that

$$\delta_S(R^v, T) = R_T^v \sum_{L \subseteq S} (-1)^{|S|-|L|} \prod_{j \in L} q_{j,v}.$$

Plugging this representation into the definition of SII with

$$\omega(t) := \frac{1}{(n - |S| + 1) \cdot \binom{n-|S|}{|T|}}$$

yields

$$I^{\text{SII}}(R^v, S) := \sum_{T \subseteq N \setminus S} \omega(|T|) \delta_S(f, T) = \left(\sum_{L \subseteq S} (-1)^{|S|-|L|} \prod_{j \in L} q_{j,v} \right) \sum_{T \subseteq N \setminus S} \omega(|T|) R_T^v.$$

This sum is further simplified to

$$\sum_{T \subseteq N \setminus S} \omega(|T|) R_T^v = R_\emptyset^v \sum_{T \subseteq N \setminus S} \omega(|T|) \prod_{j \in T} q_{j,v} = R_\emptyset^v \sum_{k=0}^{n-|S|} \omega(k) \sum_{T \subseteq N \setminus S, |T|=k} \prod_{j \in T} q_{j,v}. \quad (9)$$

We thus need to show that this term is equal to ψ applied on

$$\frac{G_v^{\text{SP}}(y)}{\prod_{j \in S} (q_{j,v} + y)} = R_\emptyset^v \frac{\prod_{j \in \mathcal{F}(R^v)} (q_{j,v} + y)}{\prod_{j \in S} (q_{j,v} + y)} = R_\emptyset^v \prod_{j \in \mathcal{F}(R^v) \setminus S} (q_{j,v} + y).$$

With $d := |\mathcal{F}(R^v)|$ using the scale invariance of the SP, we can scale this polynomial to the degree of n , i.e. multiplying it with $(1 + y)$ for every feature not present in $\mathcal{F}(R^v)$, i.e. $N \setminus \mathcal{F}(R^v)$. This yields with the above that

$$\frac{G_v^{\text{SP}}(y) \odot (1 + y)^{n-d}}{\prod_{j \in S} (q_{j,v} + y)} = R_\emptyset^v \left(\prod_{j \in \mathcal{F}(R^v) \setminus S} (q_{j,v} + y) \right) \odot (1 + y)^{n-d} = R_\emptyset^v \prod_{j \in N \setminus S} (q_{j,v} + y),$$

where the last line follows from the fact that $q_{j,v} = 1$ if the feature is not present in the path. As ψ weights the coefficients of this polynomial, we proceed by evaluating the coefficients. When writing the product as a sum, every combination $T \subseteq N \setminus S$ of features appears exactly once, and thus

$$\prod_{j \in N \setminus S} (q_{j,v} + y) = \sum_{T \subseteq N \setminus S} \left(\prod_{j \in T} q_{j,v} \right) y^{n-|T|} = \sum_{k=0}^{n-|S|} \left(\sum_{T \subseteq N \setminus S, |T|=k} \prod_{j \in T} q_{j,v} \right) y^{n-|S|-k} = \sum_{k'=0}^{n-|S|} \left(\sum_{T \subseteq N \setminus S, |T|=n-|S|-k'} \prod_{j \in T} q_{j,v} \right) y^{k'}$$

Hence, it follows for the polynomial of degree $n - |S|$ that

$$\begin{aligned}
\psi \left(\frac{G_v^{\text{SP}}(y)}{\prod_{j \in S} (q_{j,v} + y)} \right) &= \psi \left(\frac{G_v^{\text{SP}}(y) \odot (1 + y)^{n-d}}{\prod_{j \in S} (q_{j,v} + y)} \right) \\
&= \frac{R_\emptyset^v}{n - |S| + 1} \left\langle \sum_{k=0}^{n-|S|} \left(\sum_{T \subseteq N \setminus S}^{|T|=n-|S|-k} \prod_{j \in T} q_{j,v} \right) y^k, B_{n-|S|} \right\rangle \\
&= \frac{R_\emptyset^v}{n - |S| + 1} \sum_{k=0}^{n-|S|} \left(\sum_{T \subseteq N \setminus S}^{|T|=n-|S|-k} \prod_{j \in T} q_{j,v} \right) \frac{1}{\binom{n-|S|}{k}} \\
&= \frac{R_\emptyset^v}{n - |S| + 1} \sum_{k=0}^{n-|S|} \left(\sum_{T \subseteq N \setminus S}^{|T|=k} \prod_{j \in T} q_{j,v} \right) \frac{1}{\binom{n-|S|}{n-|S|-k}} \\
&= \frac{R_\emptyset^v}{n - |S| + 1} \sum_{k=0}^{n-|S|} \left(\sum_{T \subseteq N \setminus S}^{|T|=k} \prod_{j \in T} q_{j,v} \right) \frac{1}{\binom{n-|S|}{k}} \\
&= R_\emptyset^v \sum_{k=0}^{n-|S|} \omega(k) \sum_{T \subseteq N \setminus S}^{|T|=k} \prod_{j \in T} q_{j,v},
\end{aligned}$$

which is equal to (9) and finishes the proof. \square

A.2 Proof of Proposition 3

For the sum of coefficients of the IP, it holds

$$\kappa(H_{S,e}^{\text{IP}}) = \sum_{L \subseteq S} (-1)^{|S|-|L|} \prod_{j \in L} p_{j,e}. \quad (10)$$

If there exists $j \in S$ with $p_{j,e} = 1$, then $\kappa(H_{S,e}^{\text{IP}}) = 0$.

Proof. We first compute the coefficients of the IP as

$$H_{S,e}^{\text{IP}}(y) = \prod_{j \in S} (p_{j,e} - y) = \sum_{k=0}^{|S|} (-1)^{|S|-k} \left(\sum_{L \subseteq S}^{|L|=k} \prod_{j \in L} p_{j,e} \right) y^{|S|-k} = \sum_{k=0}^{|S|} (-1)^k \left(\sum_{L \subseteq S}^{|L|=|S|-k} \prod_{j \in L} p_{j,e} \right) y^k.$$

Hence,

$$\kappa(H_{S,e}^{\text{IP}}) = \sum_{k=0}^{|S|} (-1)^k \sum_{L \subseteq S}^{|L|=|S|-k} \prod_{j \in L} p_{j,e} = (-1)^{|S|} \sum_{k=0}^{|S|} (-1)^k \sum_{L \subseteq S}^{|L|=k} \prod_{j \in L} p_{j,e} = \sum_{L \subseteq S} (-1)^{|S|-|L|} \prod_{j \in L} p_{j,e},$$

which finishes the first part of the proof.

Now, let $p_{j_0,e} = 1$ for some $j_0 \in S$. Then,

$$\begin{aligned}
\kappa(H_{S,e}^{\text{IP}}) &= \sum_{L \subseteq S} (-1)^{|S|-|L|} \prod_{j \in L} p_{j,e} \\
&= \sum_{L \subseteq S \setminus \{j_0\}} (-1)^{|S|-|L|} \prod_{j \in L} p_{j,e} + \sum_{L \subseteq S \setminus \{j_0\}} (-1)^{|S|-|L \cup \{j_0\}|} \prod_{j \in L \cup \{j_0\}} p_{j,e} \\
&= \sum_{L \subseteq S \setminus \{j_0\}} (-1)^{|S|-|L|} \prod_{j \in L} p_{j,e} \\
&\quad - \sum_{L \subseteq S \setminus \{j_0\}} (-1)^{|S|-|L|} \prod_{j \in L} p_{j,e} \\
&= 0,
\end{aligned}$$

which finishes the proof. \square

A.3 Proof of Proposition 4

For a decision rule R^v of a leaf node v and a subset $S \subset N$, let $P_{S,v} := \bigcup_{i \in S} P_{i,v}$ and e_S^\uparrow as the closest ancestor of e in $P_{S,v}$. The SII of R^v is then

$$I^{\text{SII}}(R^v, S) = \sum_{e \in P_{S,v}} \kappa(H_{S,e}^{\text{IP}}) \psi \left(\left\lfloor \frac{G_v \odot (y+1)^{d_e - d_v}}{H_{S,e}^{\text{QP}}} \right\rfloor \right) - \kappa(H_{S,e_S^\uparrow}^{\text{IP}}) \psi \left(\left\lfloor \frac{G_v \odot (y+1)^{d_{e_S^\uparrow} - d_v}}{H_{S,e_S^\uparrow}^{\text{QP}}} \right\rfloor \right).$$

Proof. By Proposition 2, we need to show that the right hand side is equal to

$$\left(\sum_{L \subseteq S} (-1)^{|S|-|L|} \prod_{j \in L} q_{j,v} \right) \psi \left(\left\lfloor \frac{G_v}{\prod_{j \in S} (q_{j,v} + y)} \right\rfloor \right). \quad (11)$$

We can simplify the sum on the right hand side, as all terms except the last edge $e_S^* \in P_{S,v}$ cancel out, to

$$\begin{aligned} & \sum_{e \in P_{S,v}} \kappa(H_{S,e}^{\text{IP}}) \psi \left(\left\lfloor \frac{G_v \odot (y+1)^{d_e - d_v}}{H_{S,e}^{\text{QP}}} \right\rfloor \right) - \kappa(H_{S,e_S^\uparrow}^{\text{IP}}) \psi \left(\left\lfloor \frac{G_v \odot (y+1)^{d_{e_S^\uparrow} - d_v}}{H_{S,e_S^\uparrow}^{\text{QP}}} \right\rfloor \right) \\ &= \kappa(H_{S,e_S^*}^{\text{IP}}) \psi \left(\left\lfloor \frac{G_v \odot (y+1)^{d_{e_S^*} - d_v}}{H_{S,e_S^*}^{\text{QP}}} \right\rfloor \right) \\ &= \kappa(H_{S,e_S^*}^{\text{IP}}) \psi \left(\left\lfloor \frac{G_v \odot (y+1)^{d_{e_S^*} - d_v}}{\prod_{j \in S} (p_{j,e_S^*} + y)} \right\rfloor \right) \\ &= \kappa(H_{S,e_S^*}^{\text{IP}}) \psi \left(\left\lfloor \frac{G_v}{\prod_{j \in S} (p_{j,e_S^*} + y)} \odot (y+1)^{d_{e_S^*} - d_v} \right\rfloor \right) \\ &= \kappa(H_{S,e_S^*}^{\text{IP}}) \psi \left(\left\lfloor \frac{G_v}{\prod_{j \in S} (p_{j,e_S^*} + y)} \right\rfloor \right), \end{aligned}$$

where we used the scale invariance of ψ and the definition of the QP. For the last edge e_S^* in the path $P_{S,v}$ it holds that $p_{i,e_S^*} = q_{i,v}$ for all $i \in S$ and thus the argument in ψ is equal to the argument in (11). Furthermore, for the IP, we have by Proposition 3 that

$$\kappa(H_{S,e_S^*}^{\text{IP}}) = \sum_{L \subseteq S} (-1)^{|S|-|L|} \prod_{j \in L} p_{j,e_S^*} = \sum_{L \subseteq S} (-1)^{|S|-|L|} \prod_{j \in L} q_{j,v},$$

which concludes the proof. \square

A.4 Proof of Theorem 2

For $S \subset N$ let $E_S := \bigcup_{i \in S} E_i$ be the set of edges that split on any feature in S and denote e_S^\uparrow as the closest ancestor of e in $P_{S,v}$. The SII is then computed as

$$I^{\text{SII}}(f, S) = \sum_{e \in E_S} \kappa(H_{S,e}^{\text{IP}}) \psi \left(\left\lfloor \frac{G_{h(e)}}{H_{S,e}^{\text{QP}}} \right\rfloor \right) - \kappa(H_{S,e_S^\uparrow}^{\text{IP}}) \psi \left(\left\lfloor \frac{G_{h(e)} \odot (y+1)^{d_{e_S^\uparrow} - d_e}}{H_{S,e_S^\uparrow}^{\text{QP}}} \right\rfloor \right).$$

Proof. The model f can be represented as a sum of decision rules as $f(T) = \sum_{v \in \mathcal{L}(T_f)} R_T^v$ and thus by the linearity of SII and Proposition 4

$$\begin{aligned} I^{\text{SII}}(f, S) &= \sum_{v \in \mathcal{L}(T_f)} I^{\text{SII}}(R^v, S) \\ &= \sum_{v \in \mathcal{L}(T_f)} \sum_{e \in P_{S,v}} \kappa(H_{S,e}^{\text{IP}}) \psi \left(\left\lfloor \frac{G_v \odot (y+1)^{d_e - d_v}}{H_{S,e}^{\text{QP}}} \right\rfloor \right) - \kappa(H_{S,e_S^\uparrow}^{\text{IP}}) \psi \left(\left\lfloor \frac{G_v \odot (y+1)^{d_{e_S^\uparrow} - d_v}}{H_{S,e_S^\uparrow}^{\text{QP}}} \right\rfloor \right) \\ &= \sum_{v \in \mathcal{L}(T_f)} \sum_{e \in P_{S,v}} \kappa(H_{S,e}^{\text{IP}}) \psi \left(\left\lfloor \frac{G_v \odot (y+1)^{d_e - d_v}}{H_{S,e}^{\text{QP}}} \right\rfloor \right) - \kappa(H_{S,e_S^\uparrow}^{\text{IP}}) \psi \left(\left\lfloor \frac{G_v \odot (y+1)^{(d_e - d_v) + (d_{e_S^\uparrow} - d_e)}}{H_{S,e_S^\uparrow}^{\text{QP}}} \right\rfloor \right), \end{aligned}$$

where we added $0 = d_e - d_e$ to the scaling factor. Now we have that $v \in \mathcal{L}(\mathcal{T}_f) \wedge e \in E_S \Leftrightarrow v \in \mathcal{L}(h(e))$, i.e. if v is a leaf and e and edge in its path, then v is also reachable from the head of e , $h(e)$, and vice versa. We can thus change the summation to

$$I^{\text{III}}(f, S) = \sum_{e \in E_S} \sum_{v \in \mathcal{L}(h(e))} \kappa(H_{S,e}^{\text{IP}}) \psi \left(\left\lfloor \frac{G_v \odot (y+1)^{d_e - d_v}}{H_{S,e}^{\text{QP}}} \right\rfloor \right) - \kappa(H_{S,e_S^\uparrow}^{\text{IP}}) \psi \left(\left\lfloor \frac{G_v \odot (y+1)^{(d_e - d_v) + (d_{e_S^\uparrow} - d_e)}}{H_{S,e_S^\uparrow}^{\text{QP}}} \right\rfloor \right).$$

Note that $G_v \odot (y+1)^{d_e - d_v}$ yields a polynomial of degree d_e and the degree of $H_{S,e}^{\text{QP}}$ is always equal to $|S|$. Hence, the polynomials can be summed using the additivity of ψ , which yields

$$\begin{aligned} \sum_{v \in \mathcal{L}(h(e))} \kappa(H_{S,e}^{\text{IP}}) \psi \left(\left\lfloor \frac{G_v \odot (y+1)^{d_e - d_v}}{H_{S,e}^{\text{QP}}} \right\rfloor \right) &= \kappa(H_{S,e}^{\text{IP}}) \psi \left(\sum_{v \in \mathcal{L}(h(e))} \left\lfloor \frac{G_v \odot (y+1)^{d_e - d_v}}{H_{S,e}^{\text{QP}}} \right\rfloor \right) \\ &= \kappa(H_{S,e}^{\text{IP}}) \psi \left(\left\lfloor \frac{\sum_{v \in \mathcal{L}(h(e))} G_v \odot (y+1)^{d_e - d_v}}{H_{S,e}^{\text{QP}}} \right\rfloor \right) \\ &= \kappa(H_{S,e}^{\text{IP}}) \psi \left(\left\lfloor \frac{\bigoplus_{v \in \mathcal{L}(h(e))} G_v}{H_{S,e}^{\text{QP}}} \right\rfloor \right) \\ &= \kappa(H_{S,e}^{\text{IP}}) \psi \left(\left\lfloor \frac{G_{h(e)}}{H_{S,e}^{\text{QP}}} \right\rfloor \right). \end{aligned}$$

Similarly, for the other term

$$\begin{aligned} &\sum_{v \in \mathcal{L}(h(e))} \kappa(H_{S,e_S^\uparrow}^{\text{IP}}) \psi \left(\left\lfloor \frac{G_v \odot (y+1)^{(d_e - d_v) + (d_{e_S^\uparrow} - d_e)}}{H_{S,e_S^\uparrow}^{\text{QP}}} \right\rfloor \right) \\ &= \sum_{v \in \mathcal{L}(h(e))} \kappa(H_{S,e_S^\uparrow}^{\text{IP}}) \psi \left(\left\lfloor \frac{G_v \odot (y+1)^{(d_e - d_v) + (d_{e_S^\uparrow} - d_e)}}{H_{S,e_S^\uparrow}^{\text{QP}}} \right\rfloor \right) \\ &= \kappa(H_{S,e_S^\uparrow}^{\text{IP}}) \psi \left(\sum_{v \in \mathcal{L}(h(e))} \left\lfloor \frac{G_v \odot (y+1)^{(d_e - d_v) + (d_{e_S^\uparrow} - d_e)}}{H_{S,e_S^\uparrow}^{\text{QP}}} \right\rfloor \right) \\ &= \kappa(H_{S,e_S^\uparrow}^{\text{IP}}) \psi \left(\left\lfloor \frac{\sum_{v \in \mathcal{L}(h(e))} G_v \odot (y+1)^{d_e - d_v}}{H_{S,e_S^\uparrow}^{\text{QP}}} \odot (y+1)^{d_{e_S^\uparrow} - d_e} \right\rfloor \right) \\ &= \kappa(H_{S,e_S^\uparrow}^{\text{IP}}) \psi \left(\left\lfloor \frac{\bigoplus_{v \in \mathcal{L}(h(e))} G_v}{H_{S,e_S^\uparrow}^{\text{QP}}} \odot (y+1)^{d_{e_S^\uparrow} - d_e} \right\rfloor \right) \\ &= \kappa(H_{S,e_S^\uparrow}^{\text{IP}}) \psi \left(\left\lfloor \frac{G_{h(e)} \odot (y+1)^{d_{e_S^\uparrow} - d_e}}{H_{S,e_S^\uparrow}^{\text{QP}}} \right\rfloor \right), \end{aligned}$$

which finishes the proof. □

B Implementations

In the following, we describe pseudo-code for TreeSHAP-IQ for SII (Section B.1) and general CII (Section B.2) and Linear TreeSHAP (Section B.3), as well as its efficient implementation using the multipoint interpolation form for all polynomials.

B.1 TreeSHAP-IQ Algorithm

The pseudo-code for TreeSHAP-IQ is outlined in Algorithm 1.

Algorithm 1: $\text{TRAVERSE_TREE}(v, C, C^{\text{IP}}, C^{\text{QP}}, a_x)$ – TreeSHAP-IQ for SII

Require: Interaction order s , tree $\mathcal{T} = (V, E)$, leaf predictions \mathcal{V}_v and empty decision rules R_\emptyset^v for all leaf nodes $v \in \mathcal{L}(\mathcal{T})$.

For all edges $e \in E$ with label i : Weights w_e and inter-path products $p^{\text{raw}}(e) := 1 / \prod_{e' \in P_{i, h(e)}} w_{e'}$

For all edges $e \in E$ with label i and interaction subsets S with $i \in S$: Ancestors e_S^\uparrow .

```

1: if  $v$  is not root then
2:    $e, i \leftarrow$  edge with  $v$  as head and split feature  $i$  of its source (parent of  $v$ )
3:    $p_{i,e}(x) \leftarrow a_x(v)p^{\text{raw}}(e)$ 
4:    $C \leftarrow C \odot (y + p_{i,e}(x))$ 
5:    $H_{S,e}^{\text{IP}} \leftarrow C^{\text{IP}}[S] \odot (-y + p_{i,e}(x))$ , if  $i \in S$ , else  $C^{\text{IP}}[S]$ 
6:    $H_{S,e}^{\text{QP}} \leftarrow C^{\text{QP}}[S] \odot (y + p_{i,e}(x))$ , if  $i \in S$ , else  $C^{\text{QP}}[S]$ 
7:   if  $e$  has ancestor  $e_S^\uparrow$  then
8:      $p_{i,e_S^\uparrow}(x) \leftarrow a_x(h(e_S^\uparrow))p^{\text{raw}}(e_S^\uparrow)$ 
9:      $C \leftarrow \lfloor C / (y + p_{i,e_S^\uparrow}(x)) \rfloor$ 
10:     $H_{S,e}^{\text{IP}} \leftarrow \lfloor H_{S,e}^{\text{IP}} / (-y + p_{i,e_S^\uparrow}(x)) \rfloor$ ,  $\forall S : i \in S$ 
11:     $H_{S,e}^{\text{QP}} \leftarrow \lfloor H_{S,e}^{\text{QP}} / (y + p_{i,e_S^\uparrow}(x)) \rfloor$ ,  $\forall S : i \in S$ 
12:  end if
13: end if
14: if  $v$  is leaf then
15:    $G_v \leftarrow C \cdot R_\emptyset^v$ 
16: else
17:    $v_c \leftarrow$  child nodes  $v_c$  of  $v$ 
18:    $\forall v_c : a_x(v_c) \leftarrow \mathbf{1}(x \in \pi_{v_c})$ 
19:    $\forall v_c : G_{v_c} \leftarrow \text{TRAVERSE\_TREE}(v_c, C, H_{\cdot,e}^{\text{IP}}, H_{\cdot,e}^{\text{QP}}, a_x)$ 
20:    $G_v \leftarrow \bigoplus_{v_c} G_{v_c}$ 
21:   for all  $S$  with  $i \in S$  do
22:      $I[S] \leftarrow I[S] + \kappa(H_{S,e}^{\text{IP}})\psi(\lfloor G_v / H_{S,e}^{\text{QP}} \rfloor)$ 
23:     if  $e_S^\uparrow \neq \perp$  then
24:        $I[S] \leftarrow I[S] - \kappa(H_{S,e_S^\uparrow}^{\text{IP}})\psi(\lfloor (G_v \odot (1 + y)^{d_{e_S^\uparrow} - d_e}) / H_{S,e_S^\uparrow}^{\text{QP}} \rfloor)$ 
25:     end if
26:   end for
27: end if
28: return  $G_v$ 

```

B.2 Implementation of TreeSHAP-IQ for general CII

The pseudo-code of TreeSHAP-IQ applied to an arbitrary CII is outlined in Algorithm 2.

B.3 Implementation of Linear TreeSHAP

The pseudo-code of Linear TreeSHAP is outlined in Algorithm 3.

B.4 Efficient Implementation of Polynomial Operations using Multipoint Interpolation

Linear TreeSHAP and TreeSHAP-IQ rely on multiplication and division for polynomials, as well as evaluating the polynomial coefficients using ψ and κ . These operations can be efficiently implemented by storing the polynomial in a multipoint interpolation form. Instead of storing the polynomial, we store its evaluation at base points \mathcal{Y} . Multiplication and division then transfer to vector multiplication and division. The inner product of the coefficients of the polynomial, as required for ψ and κ , can be efficiently computed using the following lemma.

Algorithm 2: TRAVERSE TREE($v, C, C^{\text{IP}}, C^{\text{QP}}, a_x$) – TreeSHAP-IQ for CIIs

Require: Interaction order s , tree $\mathcal{T} = (V, E)$, leaf predictions \mathcal{V}_v and empty decision rules R_\emptyset^v for all leaf nodes $v \in \mathcal{L}(\mathcal{T})$.

For all edges $e \in E$ with label i : Weights w_e and inter-path products $p^{\text{raw}}(e) := 1 / \prod_{e' \in P_{i, h(e)}} w_{e'}$

For all edges $e \in E$ with label i and interaction subsets S with $i \in S$: Ancestors e_S^\uparrow .

- 1: **if** v is not root **then**
- 2: $e, i \leftarrow$ edge with v as head and split feature i of its source (parent of v)
- 3: $p_{i,e}(x) \leftarrow a_x(v)p^{\text{raw}}(e)$
- 4: $C \leftarrow C \odot (y + p_{i,e}(x))$
- 5: $H_{S,e}^{\text{IP}} \leftarrow C^{\text{IP}}[S] \odot (-y + p_{i,e}(x))$, if $i \in S$, else $C^{\text{IP}}[S]$
- 6: $H_{S,e}^{\text{QP}} \leftarrow C^{\text{QP}}[S] \odot (y + p_{i,e}(x))$, if $i \in S$, else $C^{\text{QP}}[S]$
- 7: **if** e has ancestor e^\uparrow **then**
- 8: $p_{i,e^\uparrow}(x) \leftarrow a_x(h(e^\uparrow))p^{\text{raw}}(e^\uparrow)$
- 9: $C \leftarrow \lfloor C / (y + p_{i,e^\uparrow}(x)) \rfloor$
- 10: $H_{S,e}^{\text{IP}} \leftarrow \lfloor H_{S,e}^{\text{IP}} / (-y + p_{i,e^\uparrow}(x)) \rfloor$, $\forall S : i \in S$
- 11: $H_{S,e}^{\text{QP}} \leftarrow \lfloor H_{S,e}^{\text{QP}} / (y + p_{i,e^\uparrow}(x)) \rfloor$, $\forall S : i \in S$
- 12: **end if**
- 13: **end if**
- 14: **if** v is leaf **then**
- 15: $G_v \leftarrow C \cdot R_\emptyset^v$
- 16: **else**
- 17: $v_c \leftarrow$ child nodes v_c of v
- 18: $\forall v_c : a_x(v_c) \leftarrow \mathbf{1}(x \in \pi_{v_c})$
- 19: $\forall v_c : G_{v_c} \leftarrow \text{TRAVERSE TREE}(v_c, C, H_{\cdot,e}^{\text{IP}}, H_{\cdot,e}^{\text{QP}}, a_x)$
- 20: $G_v \leftarrow \bigoplus_{v_c} G_{v_c}$
- 21: **for** all S with $i \in S$ **do**
- 22: $I[S] \leftarrow I[S] + \kappa(H_{S,e}^{\text{IP}})\psi(\lfloor (G_v \odot (1+y)^{n-d_e}) / H_{S,e}^{\text{QP}} \rfloor)$
- 23: **if** $e_S^\uparrow \neq \perp$ **then**
- 24: $I[S] \leftarrow I[S] - \kappa(H_{S,e_S^\uparrow}^{\text{IP}})\psi(\lfloor (G_v \odot (1+y)^{n-d_e}) / H_{S,e_S^\uparrow}^{\text{QP}} \rfloor)$
- 25: **end if**
- 26: **end for**
- 27: **end if**
- 28: **return** G_v

Algorithm 3: TRAVERSE TREE(v, C, a_x) – Linear TreeSHAP

Require: Interaction order s , tree $\mathcal{T} = (V, E)$, leaf predictions \mathcal{V}_v and empty decision rules R_\emptyset^v for all leaf nodes $v \in \mathcal{L}(\mathcal{T})$.

For all edges $e \in E$ with label i : Weights w_e , inter-path products $p^{\text{raw}}(e) := 1 / \prod_{e' \in P_{i, h(e)}} w_{e'}$ and ancestors e^\uparrow .

```
1: if  $v$  is not root then
2:    $e, i \leftarrow$  edge with  $v$  as head and split feature  $i$  of its source (parent of  $v$ )
3:    $p_e(x) \leftarrow a_x(v)p^{\text{raw}}(e)$ 
4:    $C \leftarrow C \odot (y + p_e(x))$ 
5:   if  $e$  has ancestor  $e^\uparrow$  then
6:      $p_{e^\uparrow}(x) \leftarrow a_x(h(e^\uparrow))p^{\text{raw}}(e^\uparrow)$ 
7:      $C \leftarrow \lfloor C / (y + p_{i, e^\uparrow}(x)) \rfloor$ 
8:   end if
9: end if
10: if  $v$  is leaf then
11:    $G_v \leftarrow C \cdot R_\emptyset^v$ 
12: else
13:    $v_c \leftarrow$  child nodes  $v_c$  of  $v$ 
14:    $\forall v_c : a_x(v_c) \leftarrow \mathbf{1}(x \in \pi_{v_c})$ 
15:    $\forall v_c : G_{v_c} \leftarrow \text{TRAVERSE TREE}(v_c, C, a_x)$ 
16:    $G_v \leftarrow \bigoplus_{v_c} G_{v_c}$ 
17:    $\phi[i] \leftarrow \phi[i] + (p_e - 1)\psi(\lfloor G_v / (p_e + y) \rfloor)$ 
18:   if  $e_S^\uparrow \neq \perp$  then
19:      $\phi[i] \leftarrow \phi[i] - (p_e^\uparrow - 1)\psi(\lfloor G_v / (p_{e^\uparrow} + y) \rfloor)$ 
20:   end if
21: end if
22: return  $G_v$ 
```

Lemma 1 (Yu et al., 2022). *Let $p, q \in \mathbb{R}[x]_d$ with coefficients A, B , respectively. Then $\langle p, q \rangle = \langle A, B \rangle = \langle p(Y), \mathcal{V}(Y)^{-1}B \rangle$, where $\mathcal{V} \in \mathbb{R}^{d+1 \times d+1}$ corresponds to the Vandermonde matrix with entries $(\mathcal{V}(Y))_{i,j} = y_i^j$.*

As proposed by Yu et al. (2022), we use the Chebyshev points \mathcal{Y} to evaluate the polynomial, as they are optimal in terms of numerical stability. Using the Chebyshev points \mathcal{Y} , the values $\mathcal{V}(\mathcal{Y})^{-1}C$ have to be precomputed, where C refers to the coefficients of B . It is thus required to precompute as many values as the degree of the given polynomial.

C Experiments

This section contains further information on the experiments and additional results like a run-time analysis in Section C.2

C.1 Dataset and Model Descriptions

This section contains detailed information about the datasets, the required pre-processing steps and the model fitting.

Models The following tree-based models are used in our experiments.

- **XGBoost (XBG) (Chen and Guestrin 2016)**

XGBoost is an ensemble learning algorithm based on gradient boosting. It utilizes decision trees as base learners and optimizes a user-defined loss function through an iterative process. XGBoost incorporates regularization techniques to control model complexity and improve generalization. In our experiments, we relied on the XGBoost library².

- **Gradient-Boosted Tree (GBT)**

Gradient-Boosted Trees are an ensemble learning method that combines multiple weak learners, here decision trees, in a sequential manner. It trains each tree to correct the errors made by the previous ones, effectively reducing the overall prediction error. We used the default parametrization of the GradientBoostingClassifier and GradientBoostingRegressor classes from the scikit-learn library (Pedregosa et al. 2011)

- **Random Forest (RF)**

Random Forest is an ensemble learning algorithm that constructs a collection of decision trees by using bootstrapped subsets of the training data and random feature selection. The predictions from individual trees are then aggregated to make a final prediction. This technique reduces overfitting and enhances model robustness compared with single decision trees. We used the default parametrization of the RandomForestClassifier and RandomForestRegressor classes for classification and regression tasks, respectively (Pedregosa et al. 2011).

- **Decision Tree (DT)**

A Decision Trees is a simple yet powerful model that makes predictions by recursively splitting the data based on the most informative features. Each internal node represents a decision based on a feature, and each leaf node represents a prediction. Decision Trees are prone to overfitting, but they serve as the fundamental building blocks for ensemble methods like Random Forest and Gradient-Boosted Trees. We used the default parameterization of the DecisionTreeClassifier and DecisionTreeRegressor classes in the scikit-learn library (Pedregosa et al. 2011).

Datasets The following dataset are used in our experiments. The data was either directly retrieved from the cited sources or via scikit-learn (Pedregosa et al. 2011) or openml (Feurer et al. 2020).

- **German Credit (Hofmann 1994)**

The German Credit dataset consists of credit applicants' information from a German bank, including 20 attributes such as age, employment status, credit history, and risk assessment outcomes. It contains 1,000 instances, and the primary prediction task involves classification to determine whether an applicant is a "good" or "bad" credit risk based on their attributes. The dataset was retrieved from the UCI repository³.

- **Bank (Moro, Cortez, and Laureano 2011)**

The Bank dataset originates from a Portuguese banking institution and includes customer-related attributes, marketing campaign details, and the outcome of customers subscribing to term deposits. It contains 45,211 instances, and the primary prediction task is classification, aiming to predict whether a customer will subscribe to a term deposit or not. We retrieved the dataset via openml and the identifier 1461.

- **Adult Census (Kohavi 1996)**

Also known as the "Census Income" or "Adult Income" dataset, it contains socio-demographic attributes of individuals along with their income levels. It contains 45,222 instances and 14 attributes, and the main prediction task is classification, aiming to predict whether an individual's income exceeds 50,000 per year. We retrieved the dataset via openml and the identifier 1590.

- **Bike (Fanaee-T and Gama 2014)**

This dataset involves information from a bike-sharing program, encompassing attributes like weather conditions, time, and bike rental counts. It contains 17,379 instances and 12 attributes, and the primary prediction task is regression, aiming to predict the count of bike rentals (a continuous value) based on the provided attributes. We retrieved the dataset via openml and the identifier 42712.

- **COMPAS (Angwin et al. 2016)**

The COMPAS dataset comprises criminal defendant attributes and recidivism predictions generated by a proprietary software tool. It contains 6,172 instances and 11 attributes, and the main prediction task is classification, aiming to predict whether a criminal defendant is likely to recidivate or not. We retrieved the "simplified" dataset from Kaggle⁴, which was

²<https://xgboost.readthedocs.io/en/stable/>

³<http://archive.ics.uci.edu/dataset/144/statlog+german+credit+data>

⁴<https://www.kaggle.com/datasets/danofner/compass>

Datasets	# Instances	# Features	Target	Performance (R^2 or Accuracy)			
				XGB	GBT	RF	DT
<i>German Credit</i>	1 000	20	{0, 1}	0.7542	0.7700	0.7533	0.6833
<i>Bank</i>	45 211	16	{0, 1}	0.9064	0.9031	0.9020	0.8941
<i>Adult Census</i>	45 222	14	{0, 1}	0.8715	0.8655	0.8568	0.8524
<i>Bike</i>	17 379	12	\mathbb{R}	0.9464	0.8442	0.8840	0.8796
<i>COMPAS</i>	6 172	11	{0, 1}	0.6695	0.6776	0.6646	0.6609
<i>Titanic</i>	891	9	{0, 1}	0.7761	0.8059	0.7947	0.7723
<i>California</i>	20 640	8	\mathbb{R}	0.8315	0.8317	0.7723	0.5945

Table 2: Overview of datasets used in the experiments with performance

presented in a blog post⁵.

- **Titanic (Dawson 1995)**

The Titanic dataset records passenger information from the ill-fated Titanic voyage, including features like age, gender, class, and survival status. It contains 891 instances and 9 attributes, and the core prediction task is classification, aiming to predict whether a passenger survived the Titanic disaster or not. The data was retrieved from Kaggle⁶.

- **California (Kelley Pace and Barry 1997)**

The California housing dataset encompasses housing-related attributes for various geographical regions in California. It contains 20,640 instances and 8 attributes, and the primary prediction task is regression, aiming to predict the median value of owner-occupied homes in California based on the provided attributes. The dataset was retrieved from `scikit-learn`.

Pre-processing and Model Training The following list contains all pre-processing steps for each dataset to reproduce the experiments. For further details we refer to the technical supplement containing all experiment scripts and these steps. We base most of the data transformation on `scikit-learn` (Pedregosa et al. 2011). All models were trained with a 70%, 30% training split and fixed random seeds.

- **German Credit:** We transform the categorical columns (“checkingstatus”, “history”, “purpose”, “savings”, “employ”, “status”, “others”, “property”, “otherplans”, “housing”, “job”, “tele”, and “foreign”) into integer values via an ordinal encoding. We binarize the label into 0 and 1.
- **Bank:** We transform the categorical columns (“job”, “marital”, “education”, “default”, “housing”, “loan”, “contact”, “day”, “month”, “campaign”, and “poutcome”) into integer values via an ordinal encoding. We binarize the label into 0 and 1. Lastly, we drop rows containing null values.
- **Adult Census:** We drop rows containing null values and transform the categorical columns (“workclass”, “education”, “marital-status”, “occupation”, “relationship”, “race”, “sex”, “native-country”, and “education-num”) into integer values.
- **Bike:** We transform the categorical columns (“season”, “year”, “month”, “holiday”, “weekday”, “workingday”, and “weather”) into integer values via an ordinal encoding and drop rows containing null values.
- **COMPAS:** The dataset was used as-is and no data transformation was applied.
- **Titanic:** We transform the categorical columns (“Sex”, “Ticket”, “Cabin”, and “Embarked”) into integer values via an ordinal encoding. We impute missing values in categorical and numerical features with the mode and median values respectively.
- **California:** The dataset was used as-is and no data transformation was applied.

⁵<https://blog.fastforwardlabs.com/2017/03/09/fairml-auditing-black-box-predictive-models.html>

⁶<https://www.kaggle.com/c/titanic/data>

C.2 Run-time Analysis

In this section, we provide a run-time analysis that validates our theoretical finding about the runtime complexity of our algorithm. As described in Section 3.1, the run-time complexity of TreeSHAP-IQ for all interactions of order s is $\mathcal{O}\left(m \cdot \ell_{\mathcal{T}} \cdot d_{\max} \cdot \binom{n-1}{s-1}\right)$, where m corresponds to the number of explanation points, d_{\max} corresponds to the maximum depth of the tree and $\ell_{\mathcal{T}}$ is the number of leaves in the tree. Clearly, there are three components that affect the run-time. First, the number of explanation points scales linearly, which is clear and will not be further considered. We thus will keep $m = 1$ in the following. Second, the *tree complexity*, given by both, the number of leaves and depth of the tree affect the complexity jointly. Note, that these values are highly dependent on each other. Third, the order of interactions affects the complexity in an exponential manner, given by the binomial coefficient. In the following, to account for irregularities in the processing times, we run every explanation 10 times and average over the run-times and show the corresponding standard deviations.

Naive Comparison We compare TreeSHAP-IQ’s run-time for computing SII values up to order $s_0 \leq 5$ with a naive computation of SII. For this comparison, we create a separate synthetic classification dataset of 5 000 samples and n number of features for all $n \in [9, 15]$ (we use the `make_classification` function from `sklearn`). We fix the tree-depth to 8 and fit a decision tree for each dataset. We make sure that the trees all consists of approximately the same number of nodes (i.e. all trees reach the maximum depth). We then compute the SII values up to order s_0 with TreeSHAP-IQ and through a naive brute force computation over all combinations of subsets. We plot the log run-time of TreeSHAP-IQ and the naive SII computation in Figure 6. The comparison shows that the naive calculation scales exponentially with the number of features, while TreeSHAP-IQ scales polynomially for each interaction order s_0 .

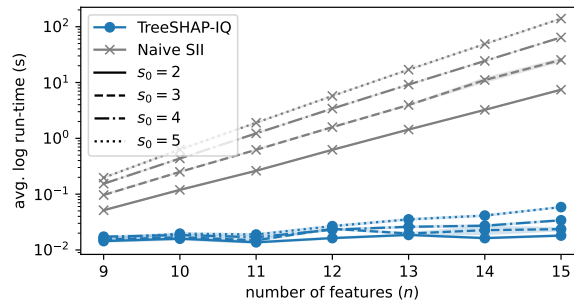


Figure 6: Log run-time of computing SII scores with TreeSHAP-IQ (blue) compared to naive calculation (grey). Naive computation scales exponentially.

Run-time by Tree Complexity We now illustrate the run-time by the tree complexity, where we compute always pairwise interactions ($s = 2$). The results are shown in Figure 7. We observe a linear relationship of the run-time and the number of vertices, as well as the number of leaves. The run-time compared with the maximum depth of the tree admits a sub-linear behavior at higher depths. This can be explained by the increasing number of paths in the DT that are shorter than the maximum depth, which increasingly occurs at higher depths. Lastly, we also observe this sub-linear behavior in terms of the number of leaves times the depth, which again is explained by the overestimation of operations.

Run-Time by Number of Interactions We illustrate the run-time depending on the number of interactions for a fixed DT of depth 20. The results are shown in Figure 8. Note that interactions are only updated, if all features have appeared in the observed path. This results in far less evaluations, especially for higher order interactions. Again, this is a consequence of the structure of the DT, where only very few paths admit the maximum depth.

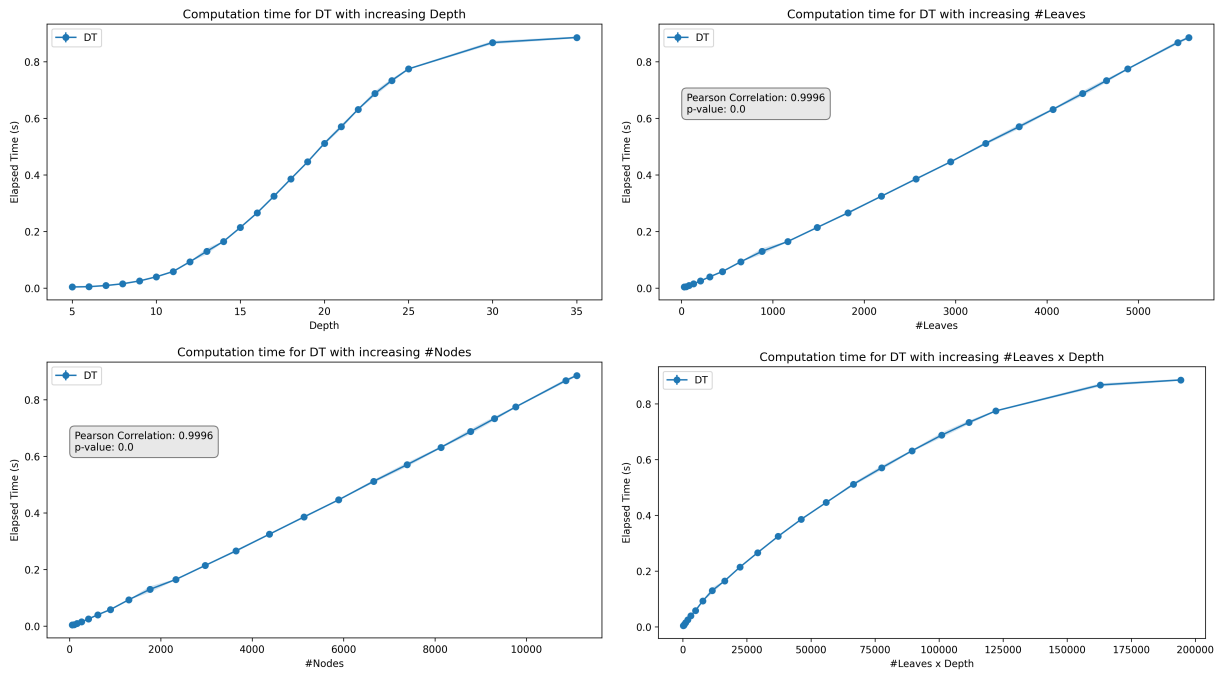


Figure 7: Run-time analysis of TreeSHAP-IQ for a single DT with varying tree complexity parameters

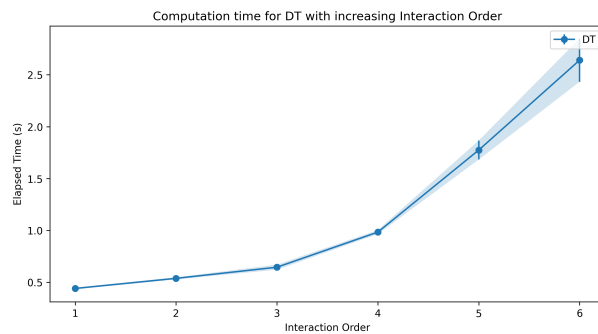


Figure 8: Run-time analysis of TreeSHAP-IQ for a single DT of depth 20 with varying interaction order

C.3 Additional n-SII Plots

In this section, we compare the strength and nature of interactions present in different model architectures. We generate the n-SII plots of order $s_0 = 6$ for the predictions of two randomly selected instances in the *Bike* and *Adult Census* dataset. The results for the *bike* dataset are shown in Figure 9 and Figure 10. The results for the *Adult Census* dataset are shown in Figure 11 and Figure 12.

We observe that the levels of interaction effects differ significantly among different predictions and different model architectures. For the *Bike* dataset, the prediction in Figure 10 has higher interaction effects than the prediction in Figure 9. Further, it can be seen that the well-performing XGB model exhibits a high amount of interaction, whereas the relatively poor-performing GBT has little interactions present. Furthermore, as expected the DT exhibits a high level of interaction effects among both predictions, as this method is not an ensemble of weak learning algorithms, which are expected to have less interaction effects.

In the *Adult Census* dataset, we observe a similar pattern for XGB and GBT. Notably, for these instances, XGB exhibits less interaction effects than RF, although the overall performance of XGB is superior. Again, we observe a high amount of higher order interactions present in the DT. However, its performance is worse than for all ensemble methods. Notably, the RF and the DT exhibit more similar interaction effects than the other models. This could be seen as an indication that the learned functional relationship in gradient boosting approaches differs from classical DT learning schemes. However, observing two local interaction effects does not allow to conclude this rigorously, which would require a global study of interaction effects in the corresponding models.

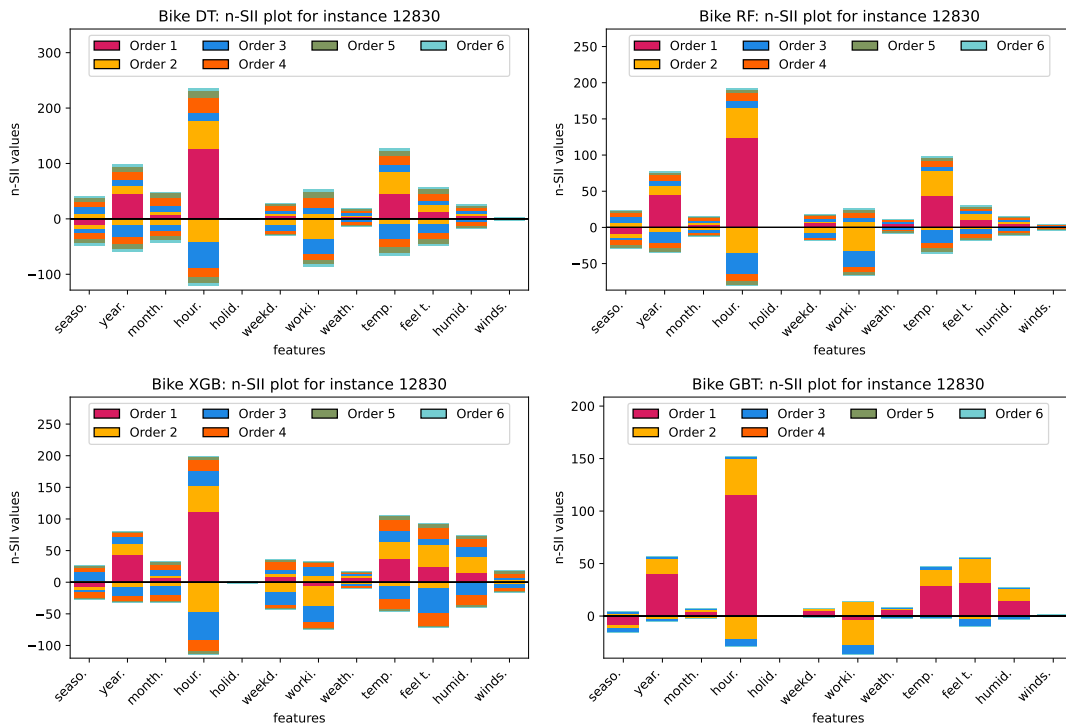


Figure 9: n-SII values up to order $s_0 = 6$ for a randomly selected instance of the *Bike* dataset.

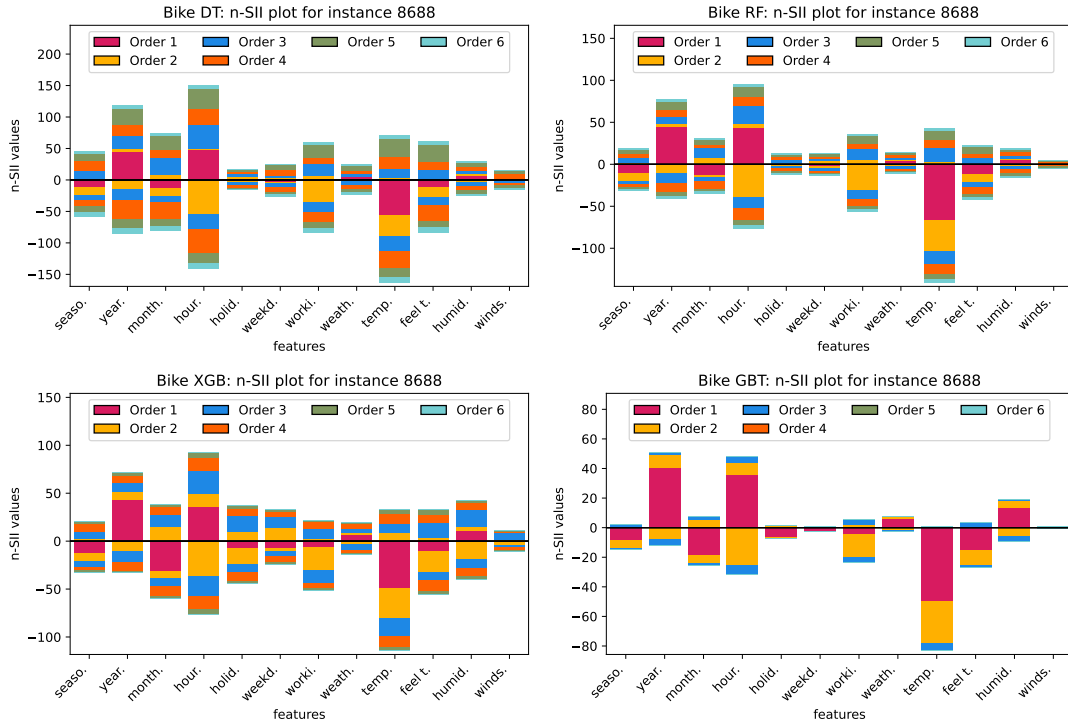


Figure 10: n-SII values up to order $s_0 = 6$ for a randomly selected instance of the *Bike* dataset.

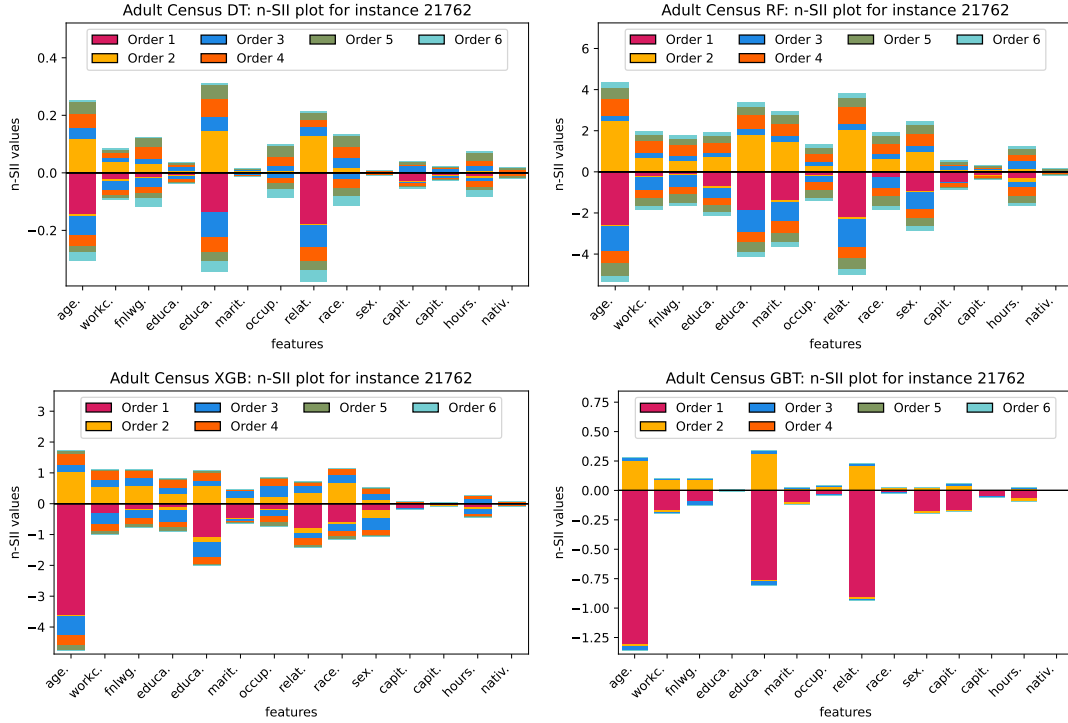


Figure 11: n-SII values up to order $s_0 = 6$ for a randomly selected instance of the *Adult Census* dataset.

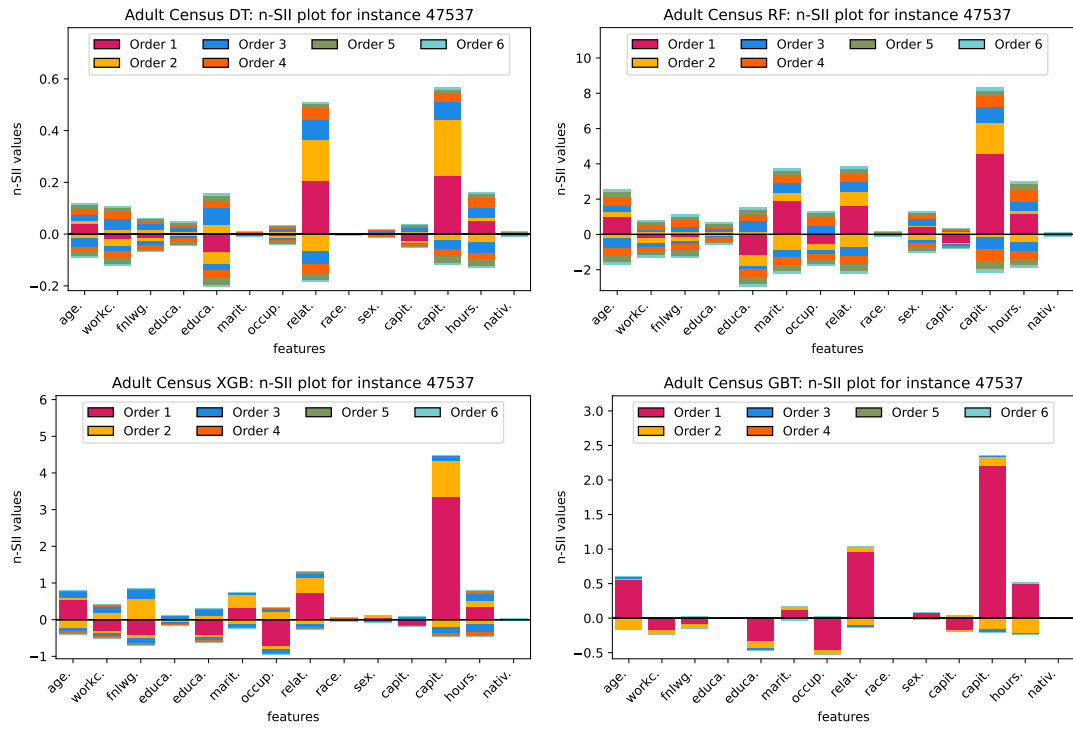


Figure 12: n-SII values up to order $s_0 = 6$ for a randomly selected instance of the *Adult Census* dataset.

C.4 Further Experimental Results on Benchmark Datasets

German Credit Dataset We display n-SII interaction effects up to order $s_0 = 2$ in a network plot and effects up to order $s_0 = 3$ in a waterfall chart for two randomly selected instances of the *German Credit* dataset predicted with a XGB. The results are shown in Figure 13.

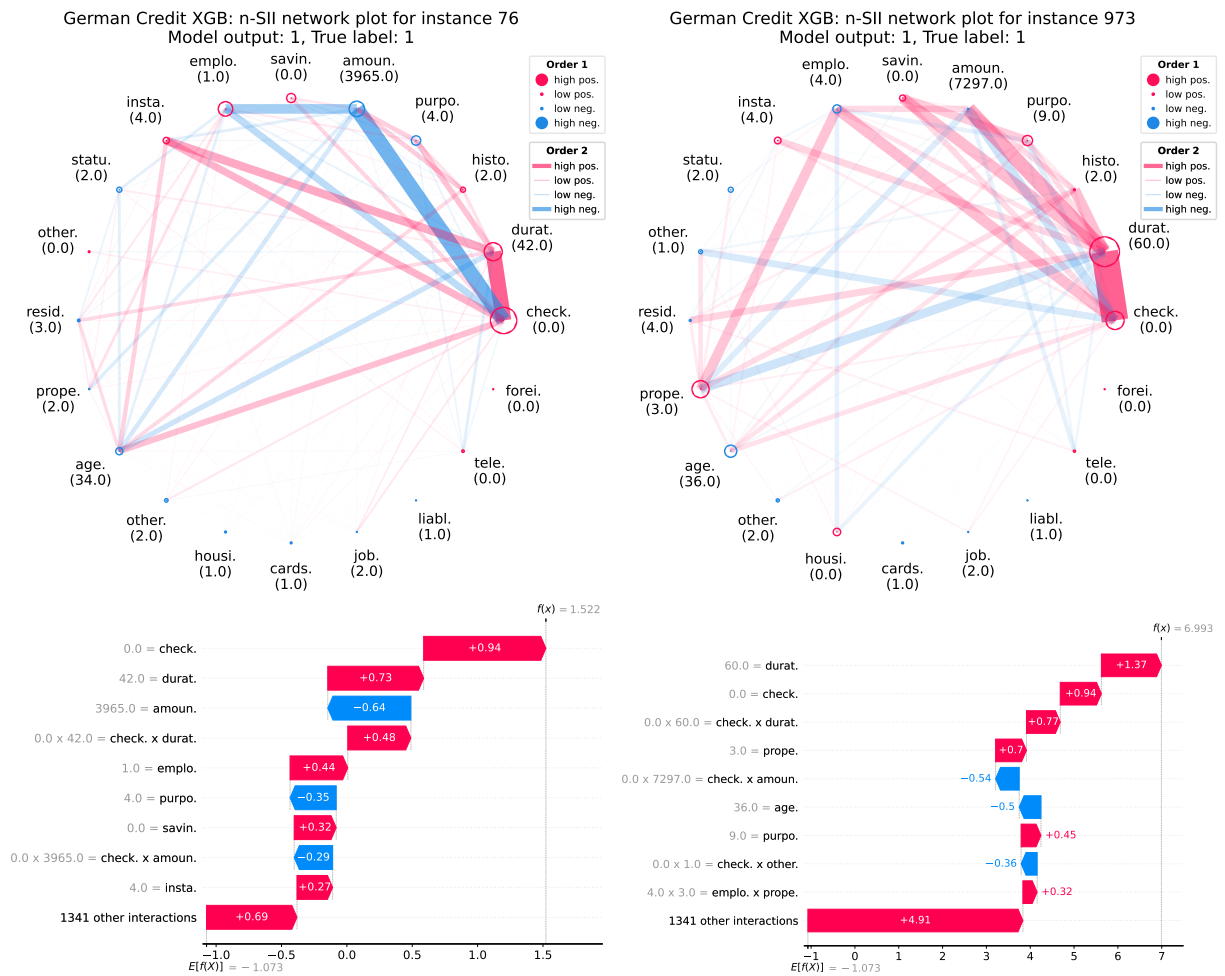


Figure 13: n-SII scores for order $s_0 = 2$ in a network plot (top) and $s_0 = 3$ in a waterfall chart (bottom) for two randomly selected instances of the *German Credit* dataset predicted with a XGB

Bank Dataset We display n-SII interaction effects up to order $s_0 = 2$ in a network plot and effects up to order $s_0 = 3$ in a waterfall chart for two randomly selected instances of the *Bank* dataset predicted with XGB. The results are shown in Figure 14.

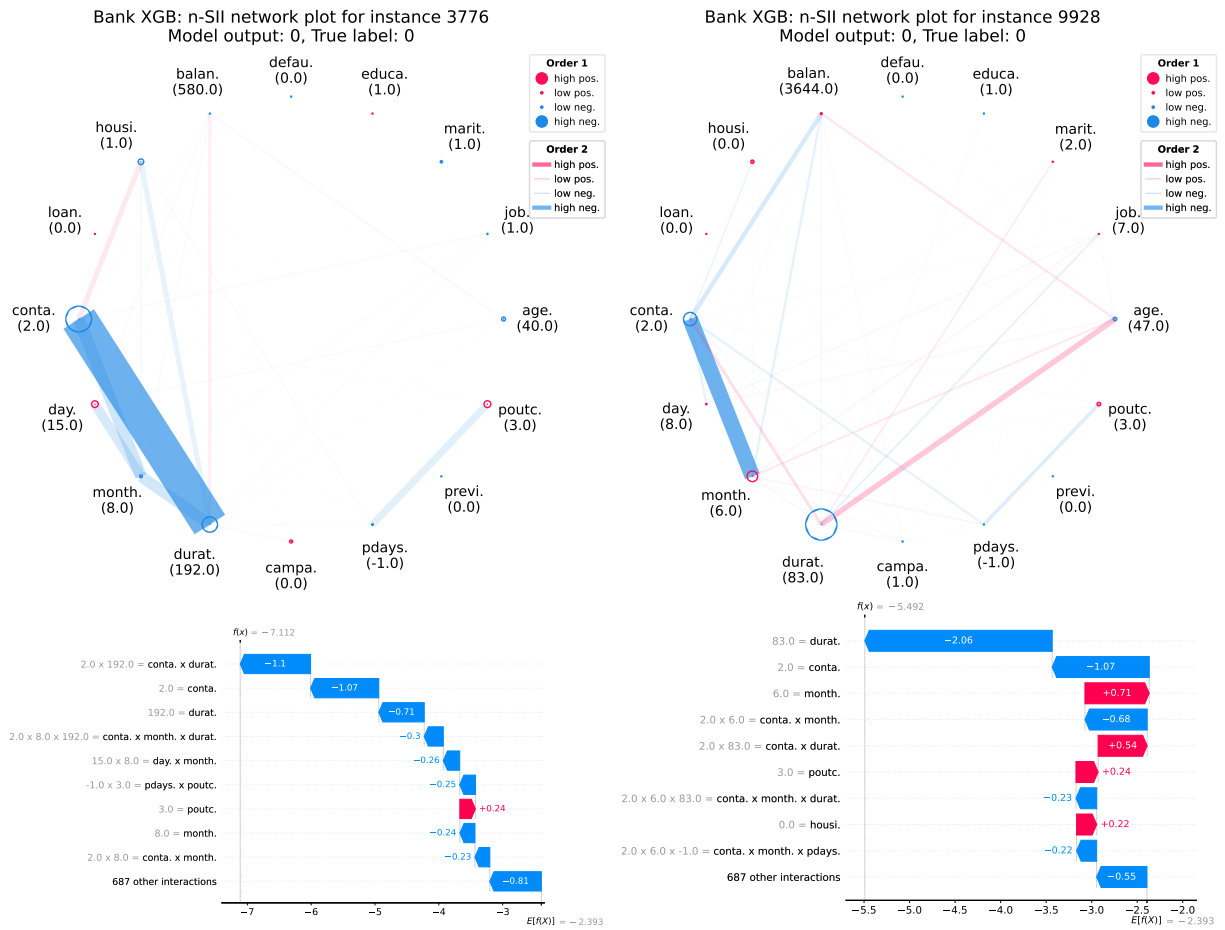


Figure 14: n-SII scores for order $s_0 = 2$ in a network plot (top) and $s_0 = 3$ in a waterfall chart (bottom) for two randomly selected instances of the *Bank* dataset predicted with XGB

Adult Census Dataset We display n-SII interaction effects up to order $s_0 = 2$ in a network plot and effects up to order $s_0 = 3$ in a waterfall chart for two randomly selected instances of the *Adult Census* dataset predicted with XGB. The results are shown in Figure 15.

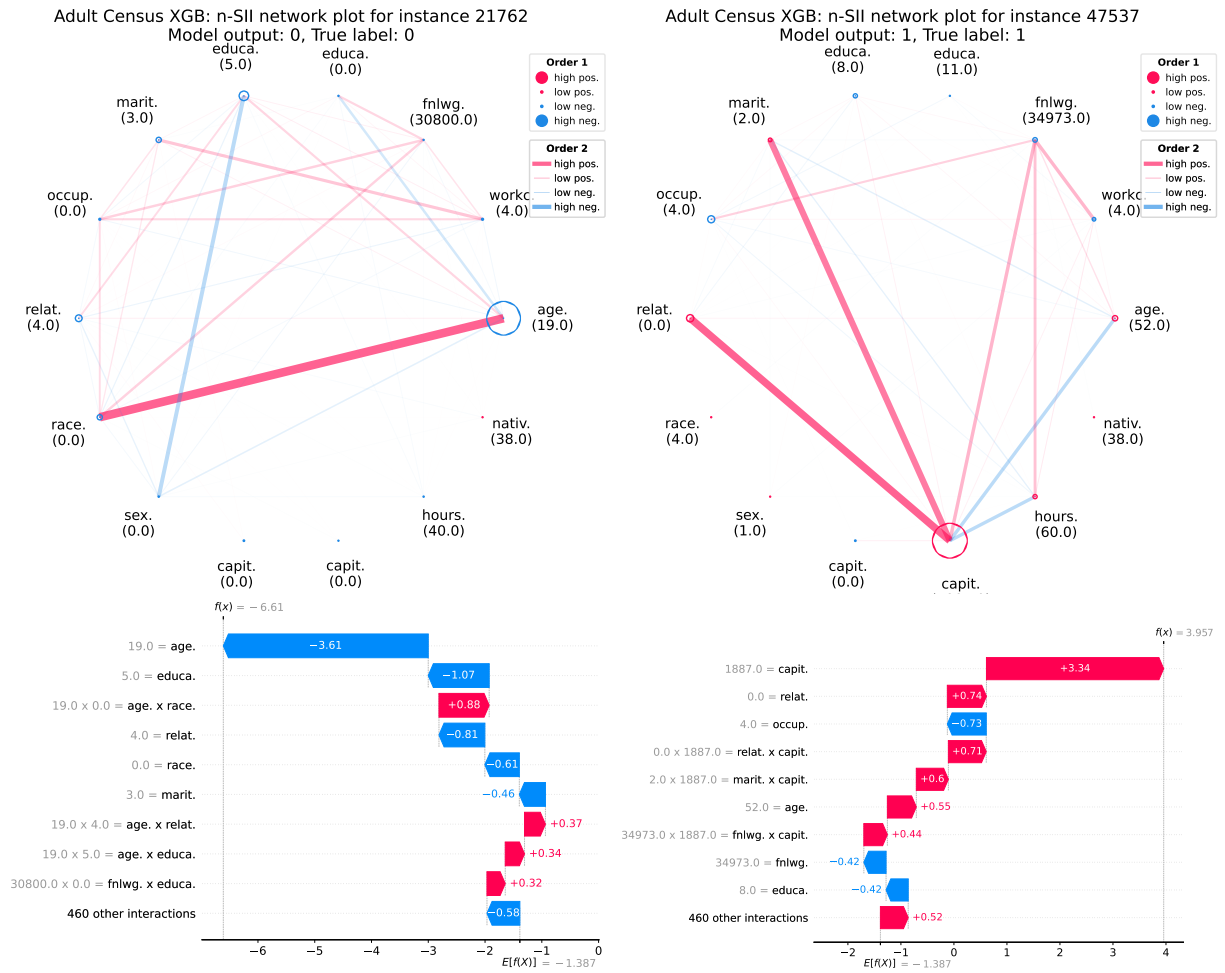


Figure 15: n-SII scores for order $s_0 = 2$ in a network plot (top) and $s_0 = 3$ in a waterfall chart (bottom) for two randomly selected instances of the *Adult Census* dataset predicted with XGB

Bike Dataset We display n-SII interaction effects up to order $s_0 = 2$ in a network plot and effects up to order $s_0 = 3$ in a waterfall chart for two randomly selected instances of the *Bike* dataset predicted with XGB. The results are shown in Figure 16.

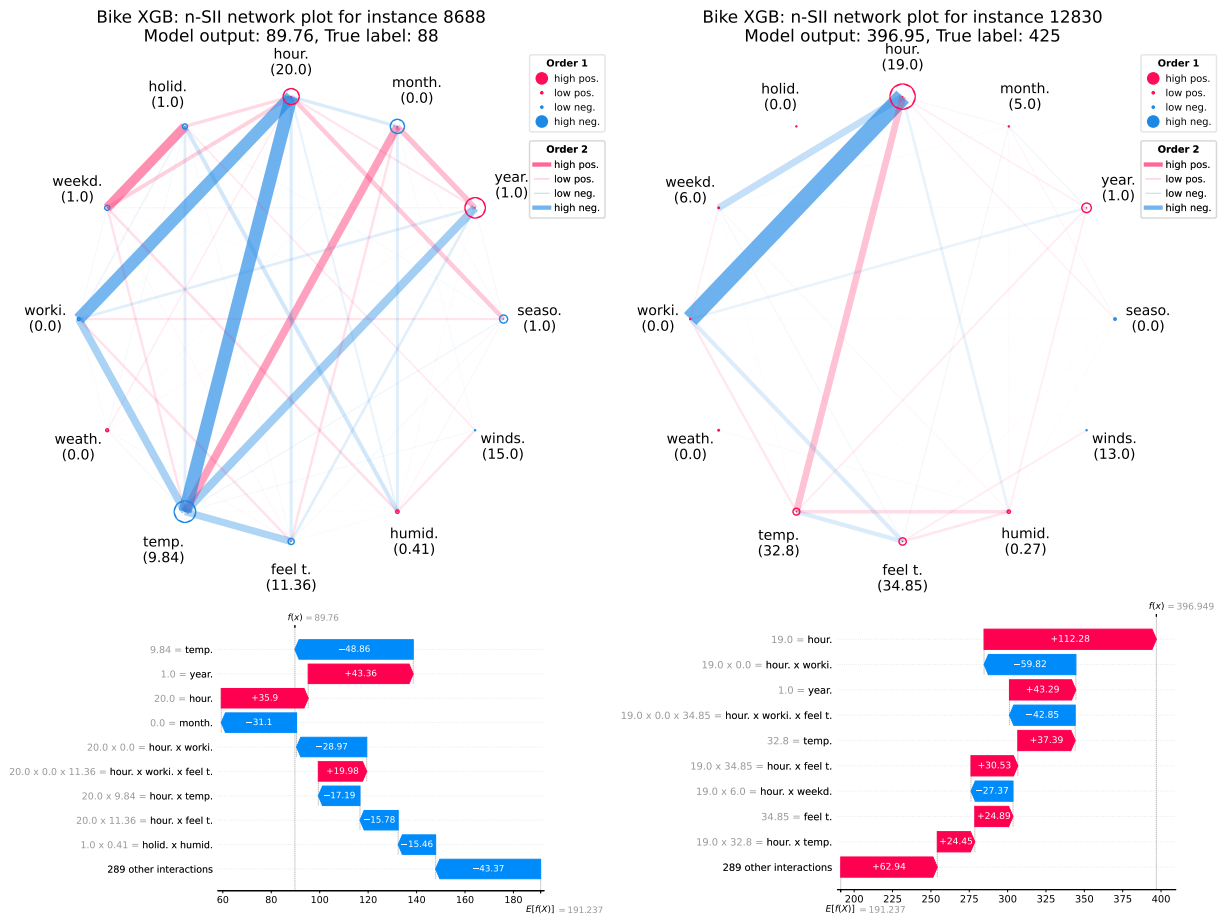


Figure 16: n-SII scores for order $s_0 = 2$ in a network plot (top) and $s_0 = 3$ in a waterfall chart (bottom) for two randomly selected instances of the *Bike* dataset predicted with XGB

COMPAS Dataset We display n-SII interaction effects up to order $s_0 = 2$ in a network plot and effects up to order $s_0 = 3$ in a waterfall chart for two randomly selected instances of the *COMPAS* dataset predicted with a GBT. The results are shown in Figure 17.

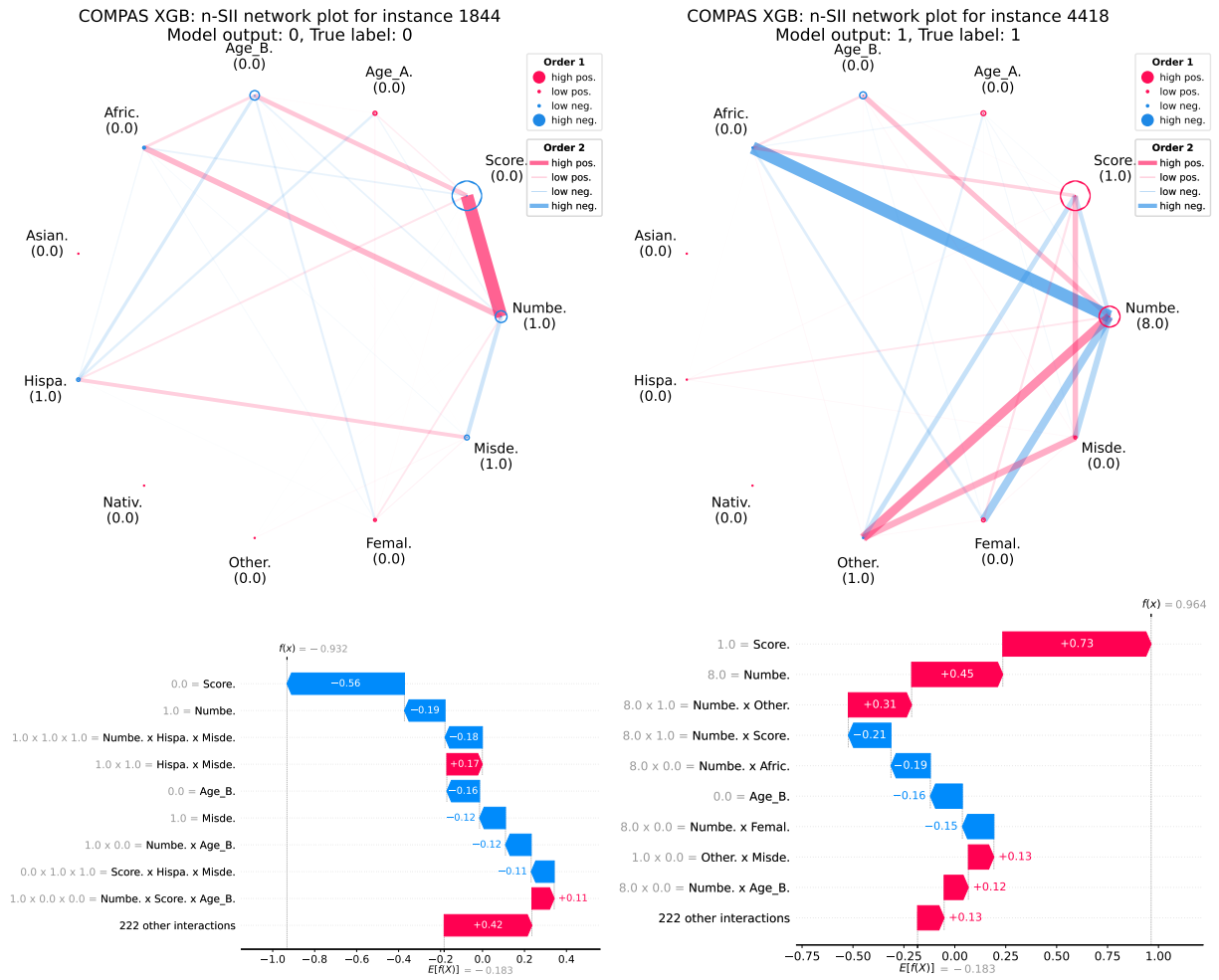


Figure 17: n-SII scores for order $s_0 = 2$ in a network plot (top) and $s_0 = 3$ in a waterfall chart (bottom) for two randomly selected instances of the *COMPAS* dataset predicted with a GBT

Titanic Dataset We display n-SII interaction effects up to order $s_0 = 2$ in a network plot and effects up to order $s_0 = 3$ in a waterfall chart for two randomly selected instances of the *Titanic* dataset predicted with a DT. The results are shown in Figure 18.

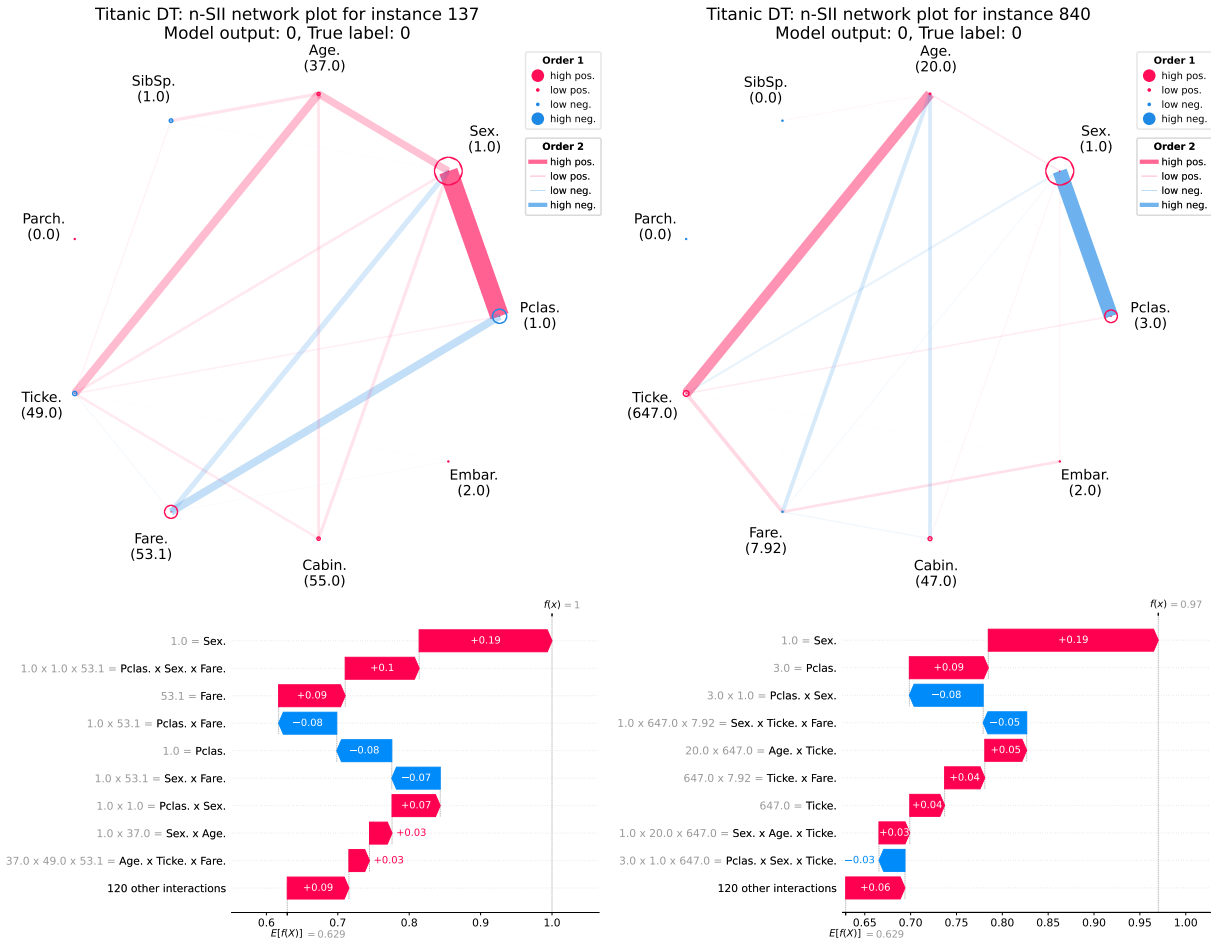


Figure 18: n-SII scores for order $s_0 = 2$ in a network plot (top) and $s_0 = 3$ in a waterfall chart (bottom) for two randomly selected instances of the *Titanic* dataset predicted with a DT

California We display n-SII interaction effects up to order $s_0 = 2$ in a network plot and effects up to order $s_0 = 3$ in a waterfall chart for two randomly selected instances of the *California* dataset predicted with a GBT. The results are shown in Figure 19.

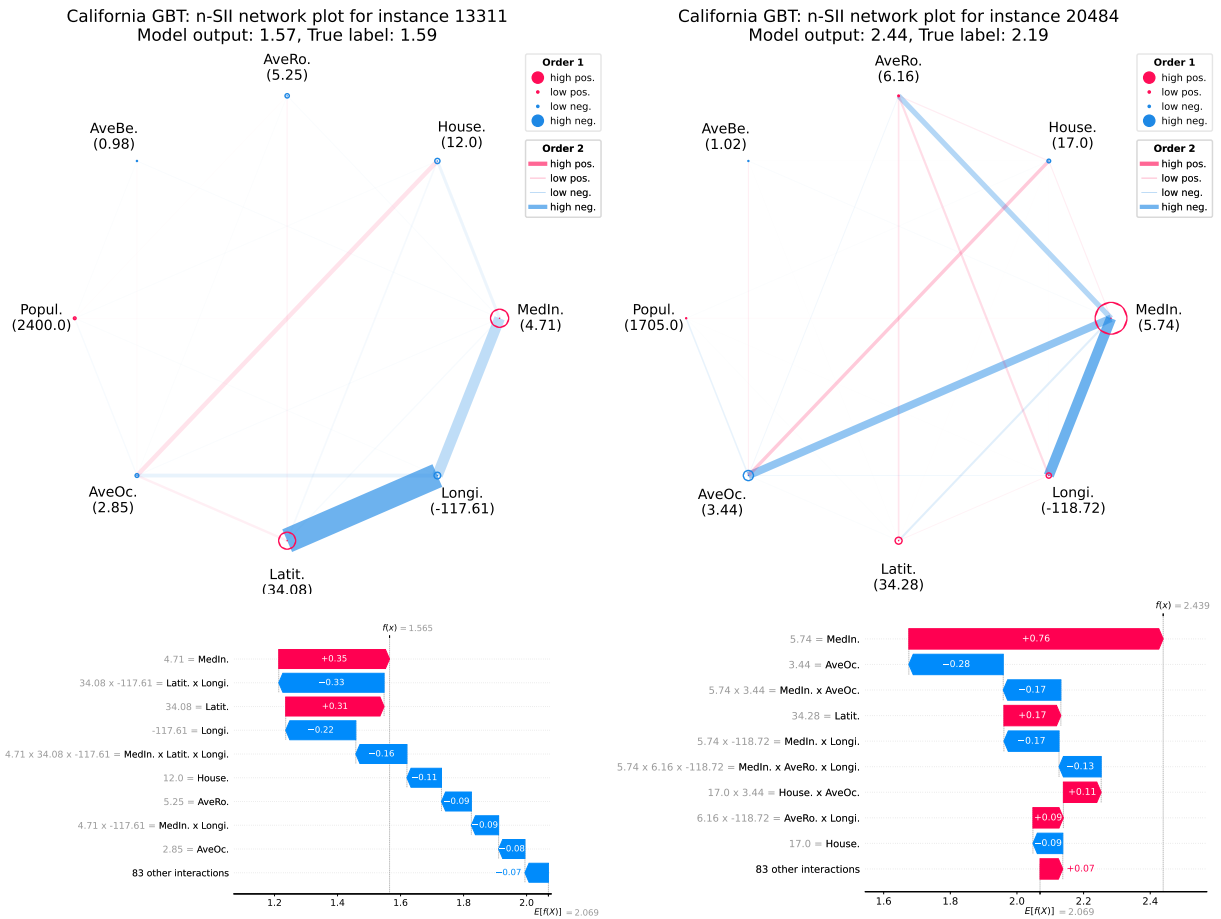


Figure 19: n-SII scores for order $s_0 = 2$ in a network plot (top) and $s_0 = 3$ in a waterfall chart (bottom) for two randomly selected instances of the *California* dataset predicted with a GBT