



# XorSHAP: Privacy-Preserving Explainable AI for Decision Tree Models

Dimitar Jetchev<sup>1</sup> and Marius Vuille<sup>2</sup>

<sup>1</sup> Input Output Global (IOG), Lausanne, Switzerland

<sup>2</sup> Arcium, Zug, Switzerland

**Abstract.** *Explainable AI (XAI)* refers to the development of AI systems and machine learning models in a way that humans can understand, interpret and trust the predictions, decisions and outputs of these models. A common approach to explainability is *feature importance*, that is, determining which input features of the model have the most significant impact on the model prediction. Two major techniques for computing feature importance are LIME (Local Interpretable Model-agnostic Explanations) and SHAP (SHapley Additive exPlanations). While very generic, these methods are computationally expensive even when the data is not encrypted. Applying them in the privacy-preserving setting when part or all of the input data is private is therefore a major computational challenge. In this paper, we present XorSHAP - the first practical data-oblivious algorithm for computing SHAP values for decision tree ensemble models. The algorithm is applicable in various privacy-preserving settings such as SMPC, FHE and differential privacy. Our algorithm has complexity  $O(TMD2^D)$ , where  $T$  is the number of decision trees in the ensemble,  $D$  is the depth of the decision trees and  $\tilde{M}$  is the maximum of the number of features  $M$  and  $2^D$  (the number of leaf nodes of a tree), and scales to real-world datasets. We implement the algorithm in the semi-honest Secure Multiparty Computation (SMPC) setting with full threshold using Inpher's Manticore framework. Our implementation simultaneously computes the SHAP values for 100 samples for an ensemble of  $T = 60$  trees of depth  $D = 4$  and  $M = 100$  features in just 7.5 minutes, meaning that the SHAP values for a single prediction are computed in just 4.5 seconds for the same decision tree ensemble model. Additionally, it is parallelization-friendly, thus, enabling future work on massive hardware acceleration with GPUs.

**Keywords:** Explainable AI · model explainability · gradient boosting decision trees · SHAP values · secure multiparty computation

## 1 Introduction

*Explainable AI (XAI)* refers to the design and development of AI systems and machine learning methods that allow for human understanding and interpretation of the decisions and predictions made by these models [Mol22]. It provides insights into how a model generates its predictions, which features it considers important as well as the reasoning behind the decisions. The main challenge of explainable AI is thus to create explainable models without compromising model accuracy.

While explaining the decision may be transparent for simple models such as linear and logistic regression, the more complex the model becomes, the more difficult it is to explain its predictions.

There are several approaches to model explainability:

---

E-mail: [dimitar.jetchev@iohk.io](mailto:dimitar.jetchev@iohk.io) (Dimitar Jetchev), [marius@arcium.com](mailto:marius@arcium.com) (Marius Vuille)



- *Feature importance*: this approach determines the relative importance or contribution of each feature in the prediction of a model.
- *Local explainability*: a method that provides explanations for individual predictions, showing which features had the most significant impact on the specific instance.
- *Global explainability*: deriving interpretable rules or decision tree models that mimic complex machine learning models; even if the derived interpretable models can be very different, their objective is to explain any prediction of the black-box models.

Techniques such as LIME (Local Interpretable Model-agnostic Explanations) [RSG16] (as well as the official `github` repository<sup>1</sup>) and SHAP (SHapley Additive exPlanations) [LL17] are key for supporting some of the above approaches by generating explainability metrics at the instance level considering contributions of individual features.

Explainable AI is an important aspect of *Trustworthy AI*, a broader concept aiming at developing, building and deploying reliable, ethical, transparent and accountable AI systems. Trustworthy AI encapsulates various principles and guidelines to address important objectives such as ethical and responsible AI, transparency, explainability, fairness, bias mitigation, robustness, data privacy as well as accountability and governance.

The recently voted EU AI Act [Act23] is a regulatory framework proposed by the European Union (EU) to ensure the proper use of Trustworthy AI. It emphasizes the importance of providing explanations for AI decisions, especially in high-risk applications and recognizes that individuals should have the right to understand the logic, significance, and consequences of automated decisions that affect them. This is particularly relevant in areas such as finance, healthcare, industrials and legal domains where the impact of AI decisions can be significant.

In the context of the EU AI Act, techniques like SHAP can help fulfill the requirement of providing meaningful explanations for AI decisions, enable stakeholders to gain insights into how specific features influence the model’s predictions and make informed judgments about the fairness, bias, or potential risks associated with the AI system [EDPS23].

The computation of SHAP values can be rather expensive since it involves computing predictions with models built on every possible subset of features. This can become particularly challenging and even infeasible for models with large number of features and samples. One often uses approximation methods such as Monte Carlo or sampling subsets of features [LEL18, SK14]. The complexities can become even larger if the input data is sensitive. This occurs very often in use cases in healthcare, financial technology as well as manufacturing where the combination of explainability and privacy of the input data is particularly relevant.

Thus, the problem of designing and implementing efficient and scalable privacy-preserving algorithms for computing feature importance metrics such as SHAP values is of primary importance.

## 1.1 Our contributions

In this paper, we address the above problem by proposing a novel privacy-preserving algorithm, XorSHAP, for computing SHAP values for decision tree ensemble models. Our algorithm is inspired by the TreesHAP algorithm [LEL18] but relies on several data-independent (oblivious) operations such as additions, multiplications, divisions, comparisons as well as sorting permutations.

The algorithm is generic in the sense that it applies to a multitude of privacy-preserving settings such as Secure Multiparty Computation (SMPC), Fully Homomorphic Encryption (FHE), Differential Privacy (DP) among others. The only requirement is that the setting

<sup>1</sup><https://github.com/marcotcr/lime>

supports data-oblivious computation of a list of basic functions (primitives) - see the introduction to Section 3 for details.

Just as in `TreeSHAP`, the major complexity gain of `XorSHAP` comes from the fact that, instead of summing over all subsets of  $\{1, \dots, M\}$  (thus, introducing a factor  $2^M$  in the complexity), one sums over all  $2^D$  possible paths in the tree and then over all subsets of the set of feature indices of the split nodes in the path (a set of size at most  $D$ ). Here,  $M$  is the number of features and  $D$  is the depth of the binary decision trees in the ensemble model.

Our main contribution is in the way we make that part oblivious, thus, resulting in complexity that no longer has the large exponential factor  $2^M$  where  $M$  is the number of features, thus, resulting in total complexity of  $O(T\widetilde{M}D2^D)$ , where  $T$  is the number of trees,  $D$  is the depth of each binary decision tree in the ensemble model and  $\widetilde{M} = \max(M, 2^D)$ .

The efficiency of our privacy-preserving approach uses several key ideas and techniques: 1) efficient oblivious permutations and oblivious sorting algorithms; 2) computing the conditional proportion of paths  $W_{\ell,S}$  in data-independent manner using the set of distinct feature indices for a given path (a set of cardinality at most  $D$ ), a slightly extended set of cardinality exactly  $D$  and a special ordering on the power set of the extended set from Section 3.2.2. The latter is described precisely in Lemma 8.

We implement our algorithm in the Secure Multiparty Computation (SMPC) setting using the `Manticore` framework [BCD<sup>+</sup>23] in the semi-honest setting. Our particular implementation is thus applicable to the most generic data distribution setup where input data comes from two or more private data sources and is either horizontally split among the owners (i.e., every owner has different samples/rows sharing the same features/columns), vertically split (every owner has complete set of features for all the samples/rows) or any combination of the two. The goal is to compute SHAP values for a given decision tree ensemble model in the MPC setting using the fully stacked dataset and without revealing private information. Although the private data owners may play the role of the computing parties in some settings, this is not a strict requirement. In fact, the `Manticore` framework supports scenarios where private data owners secret share data among a set of compute parties that is not necessarily the same as the set of data owners. The currently used security model is semi-honest, that is, the players execute the exact steps of the algorithm with full-threshold security meaning that if all players except one decide to collude, the data of the non-colluding player is still protected. In Section 5, we demonstrate the efficiency and scalability of our implementation by presenting some preliminary benchmarks.

## 1.2 Security model

Since `XorSHAP` is a generic data-independent algorithm, it is suitable for different privacy-preserving settings such as FHE, SMPC, DP, etc. The security and privacy claims of a particular implementation are therefore the ones of the underlying privacy-preserving framework. For instance, in the particular implementation in SMPC described in Section 5, the security model is inherited from the `Manticore` framework and not specifically needed for our construction. Note that other implementations with SMPC may yield different security models such as active security models or interactive triplet generation.

The security model of our particular SMPC implementation is full-threshold, semi-honest with trusted dealer or interactive triplet generation in the offline phase [BCD<sup>+</sup>23]. Here, full-threshold means that security is guaranteed against a collusion of any proper subset of the computing parties. Semi-honest assumes that the computing parties are following the protocol and executing their assigned computations correctly; yet, they are curious in the sense that they may try to learn as much information as possible from the information they receive during the execution of the protocol. Finally, the generation of the auxiliary random masking data can be done via a trusted dealer (a party that is not

allowed to collude with any of the computing parties) or via an interactive protocol run by the computing parties.

Note that the semi-honest security model is sufficient for use cases where parties are incentivized to perform the SMPC computation correctly. Examples of such use cases are multiple banks collaborating to detect fraudulent transactions or multiple healthcare institutions collaborating on medical data for training AI models.

In a typical computation of the SHAP, one only uses the prediction function and the feature values (for the samples to be explained). XorSHAP protects either the model or the feature values (or both) and still computes the feature importance metrics, the reason why we refer to it as privacy-preserving computation of explainability metrics.

It is natural to ask how much information about the private input data is revealed if one reveals the output SHAP values. In typical privacy-preserving applications, however, the SHAP values are rarely revealed. Instead, they are obviously sorted using any privacy-preserving sorting algorithm and only the features with the highest importance metrics are revealed (without the values of the actual metrics).

### 1.3 Prior art

While initial attempts have been made to address the question of privacy-preserving explainable AI, this area of research on the intersection of machine learning and data privacy is relatively nascent and will likely grow with the expanding interest on Trustworthy and Ethical AI.

Prior work uses differential privacy in XAI and more specifically, in the computation of SHAP, in order to preserve the privacy of the input data [WAYS22]. This approach relies on a type of Monte Carlo algorithm known as Layered Shapley Algorithm based on the diminishing return property. Intuitively, the latter guarantees that small coalitions provide sufficient information on the utility of a data point and that not much remains to be gained by examining large coalitions, where the marginal contributions (Section 2.1.1) is in any way guaranteed to be small. This method is also adapted to the differential privacy setting. It will be interesting to explore how the Layered Shapley Algorithm can be adapted to other privacy-preserving techniques (e.g., SMPC using primitives similar to the ones needed for XorSHAP) whenever accuracy of the result is of importance.

Another important research direction has been in the area of federated machine learning. In [Wan19], the author proposes a method for computing feature importance via SHAP in the setting of vertical federated learning with two parties, host and guest, where the host owns part of the features, the guest owns the remaining features and the host wants to interpret a specific prediction the model makes. In order for the host to achieve the latter, the guest *à priori* needs to send the feature importance for its entire feature space to the host which may reveal too much information to the host. The idea of [Wan19] is to replace the feature space with a single unified federated feature, letting the host compute a metric for its own feature space together with that single federated feature. While this is not computing exactly the SHAP model explainability metrics, experiments show that it gives some approximate information matching the feature importance one would get by using the entire feature space of the guest. While this is an interesting approach, compared to our algorithm, it applies to very specific rather than generic data distribution settings.

In addition, the recent work [BIS<sup>+</sup>22] attempts to compute SHAP values for both horizontally or vertically partitioned data via an explainable data collaboration framework based on the model-agnostic additive feature attribution algorithm KernelSHAP in the privacy-preserving distributed machine learning setting. The starting point is that, in a horizontal federated learning scenario where each party owns part of the samples, a given party may have a biased view of the samples. Similarly, in a vertical federated learning scenario, the individual parties have a partial view of the feature space. The first proposed algorithm, **Horizontal DC-SHAP**, addresses the challenge with the discrepancies

among the contributing parties by using a sharable anchor dataset to produce consistent explanation across all collaborators. In the vertical scenario, the authors propose two different algorithms **Vertical DC-SHAP** depending on the use case: 1) feature attribution is requested at a third party for the entire set of features; 2) feature attribution is requested by one of the users for the partial set of features.

There have been some attempts to compute SHAP values in the secure multiparty computation setting in the context of data valuation and data marketplaces [TLL<sup>+</sup>22]. Recall that data valuation refers to the task of fairly compensating data owners for their contributions in the data marketplace and a suitable profit-sharing and fair scheme comes from Shapley values and cooperative game theory. The proposed secure multiparty computation approach relies on active learning, that is, a machine learning algorithm that is allowed to query the user or an external data source to label data points. While interesting as a use case and an approach, this method is not quite applicable to capture the very generic setting where no active learning is assumed.

As far as we are aware, our work is the first attempt to efficiently compute SHAP values (as feature importance metrics) for decision tree ensemble models in a generic data-oblivious manner.

## 1.4 Organization and summary of the paper

The paper is organized as follows: in Section 2 we recall some background on the general theory of Shapley values from cooperative game theory (see Section 2.1). We then specialize this theory to the case of binary decision trees in Section 2.2. We also recall some basic background on Secure Multiparty Computation with additive secret sharing in Section 2.3 as well as the **Manticore** framework [BCD<sup>+</sup>23] and its security model.

The main algorithm **XorSHAP** is presented in Section 3. This section is the heart of the paper and deserves a higher-level roadmap for better guiding the reader: our starting point is a formula for the SHAP value in terms of the proportion of paths  $W_{\ell,S}$  (conditional on a subset  $S$  of features) that flow down to leaf node  $\ell$  (Theorem 1). The remaining part of the algorithm evaluates the right-hand side of this formula in a data oblivious manner. To achieve this, we encode the different subsets of features associated with the paths ending in the leaf nodes with their characteristic vectors and use oblivious permutations to sort these. In order to control the *a priori* (unknown) Shapley weights, one introduces (see Remark 4) an auxiliary superset of size exactly  $D$  of the set distinct feature indices in a given path (a set of size at most  $D$ ) and then we use a special ordering on the power set of that auxiliary superset (Section 3.2.2) that enables us to recover the cardinality purely from the binary expansion of the corresponding index in the linear order. The data-oblivious formula for the Shapley values is then given in Theorem 2. Finally, Section 3.3 discusses the rather pedantic, but important (in the data oblivious setting) case of zero-cover split nodes. To facilitate the seemingly heavy mathematical notation, we have included a table of notation in the appendix.

We complete the paper by proper complexity analysis (presented in Section 4) as well as implementation and benchmarks on random models and input data (in Section 5). Finally, in Section 6, we mention various practical implications of **XorSHAP** on real-world use cases and describe one of these use cases in more details (a privacy-preserving root cause analysis for semiconductor manufacturing).

## 2 Background and Preliminaries

### 2.1 Background on Shapley values and SHAP

Shapley values are a concept from cooperative game theory originally introduced by Lloyd Shapley [Sha53] and used to allocate the total contribution or payoff of a group to its individual members.

#### 2.1.1 Marginal contributions.

Shapley values are formally defined in terms of expected total payoffs of subsets of players in cooperative game theory, that is, in terms of a function  $v$  that assigns to each coalition  $S$  of players the worth  $v(S)$  of that coalition. In machine learning, the analogue of players are model features and the worth function is obtained using the prediction function of a model as follows: given a subset  $S$  of features, a model  $f$  with  $M$  features as well as a sample  $x = (x_1, \dots, x_M)$ , we define the worth function

$$f_S(x) := \mathbb{E}[f(\xi) | \xi_S = x_S], \quad (1)$$

where the conditional expectation is interpreted as follows: we fix the value for feature  $j$  to  $x_j$  for all  $j \in S$  and use random values for the features  $j \notin S$ . In practice, given a training set of samples, one computes the expectation as the average over these samples.

The conditional expectation  $f_S(x)$  represents the total value of the coalition  $S$  of players. One uses that to define the marginal contribution of a player  $i \notin S$  simply as  $f_{S \cup \{i\}}(x) - f_S(x)$ .

#### 2.1.2 Definition of SHAP values.

Given the model  $f$  and the sample  $x = (x_1, \dots, x_M)$ , for each feature index  $i = 1, \dots, M$ , we define the SHAP value as

$$\phi_i(x) = \sum_{S \subset \{1, \dots, M\} \setminus \{i\}} \frac{|S|!(M-1-|S|)!}{M!} (f_{S \cup \{i\}}(x) - f_S(x)), \quad (2)$$

that is, the weighted average of all marginal contributions of feature  $i$  (the subsets  $S = \emptyset$  and  $S = \{1, \dots, M\} \setminus \{i\}$  have the largest weight, whereas subsets of size  $\lfloor (M-1)/2 \rfloor$  and  $\lceil (M-1)/2 \rceil$  respectively have the smallest weight). Computing this naïvely is infeasible as the sum is over an exponential number of subsets and for each subset, one needs to estimate the model prediction over the entire training set over which we estimate the expectation.

#### 2.1.3 Computing SHAP values in the generic case

In general (for arbitrary models), an exact computation of SHAP values is often infeasible as the definition requires summing over exponentially many coalitions of features - in our case,  $2^M$ . In practice, SHAP values are usually approximated by sampling random coalition of features and summing over the marginal contributions for these subsets.

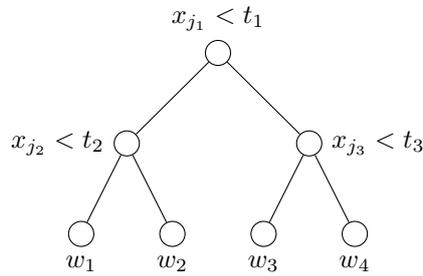
One way to approximate SHAP values is via Monte-Carlo simulations [SK14]. More specifically, this approximation method samples random feature coalitions representing random subsets of features by generating random permutations or random combinations of features and then computes the marginal contribution of each feature not included in the coalition, thus computing the chain in prediction when the feature joins the coalition. It then averages the computed marginal contributions across the sampled coalitions to provide an approximation for the SHAP value for each feature. This method has a tradeoff between the number of random coalitions sampled and the accuracy of the SHAP values.

## 2.2 SHAP values for decision trees

Even if computing SHAP values in the model-agnostic case is quite expensive, specializing to decision tree models (e.g., random forests, gradient boosted decision trees, etc.) provides some significant performance advantages which we now explain.

### 2.2.1 Binary decision trees.

Let  $M$  be a fixed positive integer (the number of features). A *binary decision tree* of depth  $D$  on the feature space  $\mathbb{R}^M$  consists of  $2^D - 1$  inner nodes (referred to as non-leaf nodes or split nodes) denoted by  $\mathcal{N}$ , and of  $2^D$  outer nodes (referred to as leaf nodes) denoted by  $L_D$  (or simply  $L$ ). Associated with each inner node  $n \in \mathcal{N}$  is a pair  $(j_n, t_n)$  of a *feature index*  $j_n \in \{1, \dots, M\}$  and a *threshold*  $t_n$  (a real number), forming the splitting criteria  $x_{j_n} < t_n$ . Associated with each leaf node  $\ell \in L$  is a *weight*  $w_\ell$  (a real number).



**Figure 1:** Binary decision tree of depth 2.

As described in [DDG<sup>+</sup>22, §2], given a sample  $x = (x_1, \dots, x_M)$  and a tree **Tree**, there is a recursive evaluation function  $\text{eval}(x, \text{Tree})$  which, for a fixed tree **Tree**, yields a piecewise constant function  $\text{Tree}: \mathbb{R}^M \rightarrow \mathbb{R}$ . Given a tree ensemble  $\{\text{Tree}^{(1)}, \dots, \text{Tree}^{(T)}\}$ , one defines the prediction on  $x$  as

$$\hat{y} = \sum_{t=1}^T \text{Tree}^{(t)}(x) \in \mathbb{R}. \quad (3)$$

### 2.2.2 SHAP values for decision trees.

Shapley values are linear in the worth function  $v$  (if two coalition games are combined, then the distributed gains correspond to the sum of the gains derived from the individual games). Moreover, the prediction of a tree ensemble model is a linear combination of the evaluation of its trees, see (3). Hence, by linearity of the mathematical expectation, it suffices to explain how to compute the SHAP values for a single tree.

We now explain how to explicitly compute the conditional expectation for a decision tree via a recurrence formula. Let **Tree** be a binary decision tree of depth  $D$  on  $M$  features, and let  $x = (x_1, \dots, x_M)$  be a sample of  $M$  feature values to be interpreted (or explained). For a subset  $S \subset \{1, \dots, M\}$  of feature indices, consider the conditional expectation  $\mathbb{E}[\text{Tree}(\xi) \mid \xi_S = x_S]$  as in (1). It can be computed as

$$\mathbb{E}[\text{Tree}(\xi) \mid \xi_S = x_S] = G_S(x; n_{\text{root}}),$$

where, for a given node  $n$  of the tree,  $G_S(x; n)$  is recursively defined as:

$$G_S(x; n) = \begin{cases} w_n & \text{if } n \text{ is a leaf node,} \\ (x_{j_n} < t_n) ? G_S(x; n_{\text{left}}) : G_S(x; n_{\text{right}}) & \text{else if } j_n \in S, \\ \frac{c_{n_{\text{left}}}}{c_n} G_S(x; n_{\text{left}}) + \frac{c_{n_{\text{right}}}}{c_n} G_S(x; n_{\text{right}}) & \text{else.} \end{cases} \quad (4)$$

Here,  $c_n$  is the *cover* of node  $n$ , i.e., the number of samples that are visiting node  $n$  upon evaluation. It is computed as the Hamming weight of the instance vector  $IV_n$  of that node (the characteristic vector of the subset of samples visiting  $n$ ).

*Remark 1.* For now we assume that all covers are positive. In the general case we must assume the possibility of some covers being zero, and we will explain in Section 3.3 how to adapt the algorithm.

Given a leaf node  $\ell \in L_D$ , we use  $\mathcal{P}_\ell = \{n_0^{(\ell)}, \dots, n_{D-1}^{(\ell)}\}$  for the (ordered) set of  $D$  split nodes in the path leading from the root  $n_0^{(\ell)}$  to leaf  $\ell$ . We also define  $\mathcal{J}_\ell = \{j_0^{(\ell)}, \dots, j_{D-1}^{(\ell)}\}$  to be the (ordered) multi-set of  $D$  feature indices for the split nodes in  $\mathcal{P}_\ell$ . Finally, let  $\mathcal{F}_\ell$  be the set of distinct feature indices for the split nodes of the path leading to leaf  $\ell$  (we thus have  $|\mathcal{F}_\ell| \leq D$ ).

**Definition 1.** We call a path  $\mathcal{P}_\ell$  *consistent* with respect to the pair  $(x, S)$  of a sample and a feature subset  $S$  if for every node  $n \in \mathcal{P}_\ell$  such that  $j_n \in S$ , the child node determined by the condition  $x_{j_n} < t_n$  coincides with the node of  $\mathcal{P}_\ell \cup \{\ell\}$  that succeeds  $n$  in the path. Otherwise, the path is called *inconsistent*.

In what follows, it is important to distinguish the following metric

$$W_{\ell,S}(x) := \begin{cases} \prod_{\substack{n \in \mathcal{P}_\ell \\ j_n \notin S}} \frac{c_{n_{\text{child}}}}{c_n} & \text{if } \mathcal{P}_\ell \text{ is consistent with } (x, S), \\ 0 & \text{else.} \end{cases} \quad (5)$$

Given a subset  $S$  and a leaf node  $\ell$ ,  $W_{\ell,S}$  computes the proportion of paths (conditional on  $x_S$ , the subset  $S$  of fixed feature values for  $S$ ) that flows down to leaf node  $\ell$ . The conditional expectation is then simply given by

$$\mathbb{E}[\text{Tree}(\xi) \mid \xi_S = x_S] = \sum_{\ell \in L} W_{\ell,S}(x) \cdot w_\ell. \quad (6)$$

Substituting (6) into the definition (2) of SHAP value yields

$$\begin{aligned} \phi_i(x) &= \sum_{S \subset \{1, \dots, M\} \setminus \{i\}} \frac{|S|!(M-1-|S|)!}{M!} \sum_{\ell \in L} (W_{\ell, S \cup \{i\}}(x) - W_{\ell,S}(x)) \cdot w_\ell \\ &= \sum_{\ell \in L} \left( \sum_{S \subset \{1, \dots, M\} \setminus \{i\}} \frac{|S|!(M-1-|S|)!}{M!} (W_{\ell, S \cup \{i\}}(x) - W_{\ell,S}(x)) \right) \cdot w_\ell. \end{aligned} \quad (7)$$

### 2.2.3 KernelSHAP and TreeSHAP

The framework SHAP for interpreting predictions was first introduced in [LL17] and was based on Shapley values. Since the latter are very expensive to compute, several approximations have been proposed.

The **KernelSHAP** method originally described in [LL17, §4] is model agnostic and is based on approximating SHAP values by sampling subsets of features containing the given feature and averaging.

In the specific case of decision trees, however, the complexity of **KernelSHAP** is  $O(T2^{D+M})$  where  $T$  is the number of trees in the ensemble,  $D$  is the tree depth and  $M$  is the total number of features.

An alternative algorithm for computing SHAP values in the case of decision trees is **TreeSHAP** [LEL18, LEC+20]. The algorithm has complexity  $O(T2^D D^2)$  and is asymptotically faster than **KernelSHAP** since in most practical applications, one assumes that the number of features  $M$  is asymptotically larger than  $D^2$ . This last assumption holds since

the depth  $D$  is often chosen to be a small number (e.g.,  $D < 10$ ), especially in the cases when the binary decision tree models are computed in the privacy-preserving settings on sensitive training data. On the other hand, the number of features in any practical scenario is often large - for instance, for financial forecasting or fraud detection models, there are hundreds or even thousands of input features.

There are further refinements to the algorithm [Yan22] resulting in **FastTreeSHAP v.1** (a 1.5x speedup and the same memory overhead as **TreeSHAP**) and **FastTreeSHAP v.2** (a 3x speedup with via some preprocessing and a slightly larger memory overhead). Furthermore, **GPUShAP** proposed in [MFH22] provides a more scalable version of **TreeSHAP** resulting in a 20x speedup.

## 2.3 Secure Multiparty Computation

Multiparty computation (MPC) is a method for cryptographic computing allowing several parties holding private data to evaluate a public function on their aggregate data while revealing only the output of the function and nothing else. Recent advances in the area make these protocols practical and suitable for real-world applications such as machine and statistical learning [BLW08, BCG<sup>+</sup>18, BCD<sup>+</sup>23, DPSZ12, EGK<sup>+</sup>20, Kel20, KOS16, KPR18, MR18, MZ17, WGC19].

In general, our **XorSHAP** algorithm is agnostic to underlying MPC method or framework. The choice of the **Manticore** MPC framework in our implementation is favorable for several reasons: 1) it provides access to Boolean arithmetic as well as arithmetic with real numbers represented using modular integers [BCD<sup>+</sup>23] or the prior floating-point numbers framework [BCG<sup>+</sup>18]; 2) the real number representation of **Manticore** via modular integers guarantees information-theoretically security; 3) **Manticore** provides oblivious sorting and oblivious permutations functionality. There are, however, other MPC libraries that support Boolean and real number arithmetic such as **SCALE-MAMBA** [aM23], **SecureML** [MZ17], **ABY** [PSSY21, MR18] as well as **SPDZ-2K** [CDE<sup>+</sup>18]. In **Manticore**, computations are split into offline and online phases. The offline phase generates random precomputed data without accessing the private data and can be performed either interactively by techniques such as oblivious transfer or by an independent party, different from the compute parties, known as trusted dealer. In the former case, the MPC protocol is slower but the security model is stronger. In the latter case, the offline phase is significantly accelerated but in order to protect the privacy of the input data, the trusted dealer is assumed not to collude with any of the computing parties. In addition, to ensure security against malicious external adversaries, all communications between the trusted dealer as well as all communication between the players during the online phase is end-to-end encrypted.

## 2.4 Secret sharing of binary decision trees

In the setting that we will apply our algorithm, there are different possibilities for the decision tree ensemble model. It can be proprietary model for a given party, or it can already be computed with a privacy-preserving algorithm such as **XorBoost** [DDG<sup>+</sup>22] or any other privacy-preserving algorithm. We refer to [DDG<sup>+</sup>22, §2.4] for details on how one can secret share a decision tree model. Note that, depending on the training algorithm, the column index  $j$  can be represented either in arithmetic shares or in Boolean shares. This is the reason why in the discussions in Section 3, we make no assumption on the secret sharing method (arithmetic or Boolean) and present a generic arithmetic-to-Boolean conversion method in Lemma 1.

Note that the actual representation of arithmetic secret shares varies depending on the SMPC implementation or framework used. For example, for most of the applications **Manticore** relies on a particular representation of real numbers with modular integers that

can differ from other schemes such as SPDZ. Our algorithms are agnostic to the particular instantiation of arithmetic secret sharing.

### 3 XorSHAP: Privacy-preserving TreeSHAP

In this section, we present the main contribution of this paper – an algorithm, XorSHAP that is a privacy-preserving variant of the TreeSHAP method [LEL18].

Our main mathematical result that yields our privacy-preserving algorithm is Theorem 2 - a data-independent formula for the SHAP values. Applying that formula in the secure multiparty computation setting requires the following operations:

- Additions of secret shared numbers (over both  $\mathbb{R}$  and  $\mathbb{Z}/2\mathbb{Z}$ ),
- Multiplications of secret shared numbers (over  $\mathbb{R}$  and  $\mathbb{F}_2$ ),
- Private divisions of secret shared numbers (over  $\mathbb{R}$ )
- Comparisons of secret shared real numbers,
- Sorting and `argsort` of secret shared real-valued vectors,
- Applying secret shared sorting permutations.

Unless otherwise stated, all subsequent variables are secret and all operations are data-independent.

#### 3.1 Preliminaries.

In order to present the algorithm, we use the notations introduced in Section 2.2.2, particularly,  $\mathcal{P}_\ell$ ,  $\mathcal{J}_\ell$  and  $\mathcal{F}_\ell$ . We also introduce some extra notation that is relevant for data-independent representation of the various structures and hence, for an SMPC-friendly algorithm.

Given a binary decision tree of depth  $D$ , recall that  $L = L_D$  denotes the set of leaf nodes of the tree. We assume that sample  $x$  is drawn from a dataset with  $M$  features. For a given leaf node  $\ell \in L$ , define

- $P_\ell$ : a binary matrix of size  $M \times D$  whose columns are the binary feature selector vectors for the split nodes of the path leading to leaf  $\ell$ . Here, a feature with index  $1 \leq j \leq M$  is encoded with the binary column vector

$$\mathbf{b}_j := \underbrace{(0, \dots, 0)_{j-1}}_{j-1}, \underbrace{(1, 0, \dots, 0)_{M-j}}_{M-j}.$$

Thus,  $P_\ell = \text{concat}(\mathbf{b}_j : j \in \mathcal{J}_\ell)$ . The matrix  $P_\ell$  is helpful for a secret shared representation of the set  $\mathcal{J}_\ell$ .

- $F_\ell$ : the characteristic vector of the set  $\mathcal{F}_\ell$ , i.e., a binary (column) vector of size  $M$  whose  $j$ th entry is 1 if  $j \in \mathcal{F}_\ell$  (i.e., feature  $j$  appears as a splitting feature in the path leading to leaf  $\ell$ ) and 0 otherwise.
- $\mathbf{o}_\ell = (o_n^{(\ell)} : n \in \mathcal{P}_\ell)$ : a binary (column) vector of size  $D$ , where  $o_n^{(\ell)} = 1$  if evaluating node  $n$  on the sample  $x$  yields the child node that is in the path  $\mathcal{P}_\ell \cup \{\ell\}$ ; otherwise  $o_n^{(\ell)} = 0$ . For simplicity, we also write  $o_{\ell,d} = o_{n_d}^{(\ell)}$  for the  $d$ th component of  $\mathbf{o}_\ell$ .

- $\mathbf{z}_\ell = (z_n^{(\ell)} : n \in \mathcal{P}_\ell)$ : a (column) vector of real numbers of size  $D$  corresponding to the cover ratios of the nodes in the path  $\mathcal{P}_\ell \cup \{\ell\}$ ; i.e., for node  $n$ ,

$$z_n^{(\ell)} = \frac{c_{n_{\text{child}}}}{c_n},$$

where  $n_{\text{child}}$  is the successor node of  $n$  in  $\mathcal{P}_\ell \cup \{\ell\}$ . For simplicity, we also write  $z_{\ell,d} = z_{n_d^{(\ell)}}^{(\ell)}$  for the  $d$ th component of  $\mathbf{z}_\ell$ . We compute  $z_n^{(\ell)}$  using the private division algorithm from [BCD<sup>+</sup>23].

**Lemma 1.** *We can compute Boolean secret shares for  $\mathbf{b}_j$  for a given arithmetic secret shared index  $j \in \{1, \dots, M\}$  by first applying any arithmetic-to-Boolean conversion to  $j$  and then a one-hot binary decoder circuit to the binary expansion of  $j$  (check [DH12, §8.1] for details).*

The lemma below computes the characteristic vector  $F_\ell$  in a data independent manner:

**Lemma 2.** *We can obtain  $F_\ell$  from  $P_\ell$  by OR-ing the columns, i.e., the following equality holds*

$$F_\ell = \bigvee_{j \in \mathcal{J}_\ell} \mathbf{b}_j.$$

*Remark 2.* Note that  $F_\ell = \bigvee_{j \in \mathcal{J}_\ell} \mathbf{b}_j$  by definition, however the right-hand side cannot be computed in a data-independent way since  $\mathcal{J}_\ell$  is not known.

Finally, we will need to compute the vector  $\mathbf{o}_\ell$  in a data-independent manner. This is achieved by the following:

**Lemma 3.** *We compute  $o_{\ell,d}$  as follows: let  $(j,t)$  be the feature index and threshold associated with node  $n_d^{(\ell)}$ , and let  $\beta = (x_j < t) \in \{0,1\}$ . Let  $n_{\text{child}}$  be the successor node of  $n_d^{(\ell)}$  in  $\mathcal{P}_\ell \cup \{\ell\}$  and let  $n_{\text{left}}$  and  $n_{\text{right}}$  be the left and right children nodes of  $n_d^{(\ell)}$  respectively (when looking at the whole tree). Then:*

$$o_{\ell,d} = \begin{cases} \beta & \text{if } n_{\text{child}} = n_{\text{left}}, \\ -\beta & \text{if } n_{\text{child}} = n_{\text{right}}. \end{cases}$$

**Lemma 4** (Complexity).

- There are  $2^{D-1}$  binary matrices  $P_\ell$  of size  $M \times D$ , obtained via Lemma 1. The cost is negligible. The number is  $2^{D-1}$  since sibling leaves share the same path of split nodes.
- There are  $2^{D-1}$  binary matrices  $F_\ell$  of size  $M \times 1$  (obtained by OR-ing the columns of  $P_\ell$ ). The number is  $2^{D-1}$  since sibling leaves share the same path of split nodes.
- Computing all  $2^D$   $\mathbf{o}_\ell$ 's requires a total of  $2^D - 1$  comparisons of the form  $t_n < x_{j_n}$ .
- Computing all  $2^D$   $\mathbf{z}_\ell$ 's requires one private division between two vectors of size  $2^D - 1 \times 1$  (there is a relation between the cover ratio of a left child node and the cover ratio of its right sibling node, so it suffices to compute the left ratios only).

### 3.2 XorSHAP.

Let  $x$  be a fixed sample. For what follows we will simply write  $W_{\ell,S}$  instead of  $W_{\ell,S}(x)$ . With the notation introduced in Section 3.1 we can express  $W_{\ell,S}$  as:

**Lemma 5.**

$$W_{\ell,S} = \prod_{\substack{n \in \mathcal{P}_\ell \\ j_n \in S}} o_n^{(\ell)} \prod_{\substack{n \in \mathcal{P}_\ell \\ j_n \notin S}} z_n^{(\ell)}.$$

*Proof.*  $\mathcal{P}_\ell$  is consistent with  $(x, S)$  if and only if  $\prod_{\substack{n \in \mathcal{P}_\ell \\ j_n \in S}} o_n^{(\ell)} = 1$ .  $\blacksquare$

*Remark 3.* Lemma 5 does not provide a way to obviously compute  $W_{\ell,S}$ . Instead of the costly (in a data-independent approach) evaluation of the criteria  $j_n \in S$  (respectively  $j_n \notin S$ ), our approach is based on a specific linear ordering of the subsets of  $\mathcal{F}_\ell$  (or even a slightly larger set), see Lemma 8.

The algorithm XorSHAP is based on the following theorem:

**Theorem 1.** *The SHAP values for a binary decision tree can be computed as follows:*

$$\phi_i(x) = \sum_{\substack{\ell \in L \\ i \in \mathcal{F}_\ell}} \left( \sum_{S \subset \mathcal{F}_\ell \setminus \{i\}} \frac{|S|!(|\mathcal{F}_\ell| - 1 - |S|)!}{|\mathcal{F}_\ell|!} \cdot W_{\ell,S} \right) \cdot t_i^{(\ell)} \cdot w_\ell, \quad (8)$$

where  $t_i^{(\ell)} = \prod_{\substack{n \in \mathcal{P}_\ell \\ j_n = i}} o_n^{(\ell)} / z_n^{(\ell)} - 1$ .

*Proof.* Since  $W_{\ell,S \cup \{i\}} = W_{\ell,S}$  whenever  $i \notin \mathcal{F}_\ell$ , summation in (7) simplifies to

$$\phi_i(x) = \sum_{\substack{\ell \in L \\ i \in \mathcal{F}_\ell}} \left( \sum_{S \subset \{1, \dots, M\} \setminus \{i\}} \frac{|S|!(M - 1 - |S|)!}{M!} (W_{\ell,S \cup \{i\}} - W_{\ell,S}) \right) \cdot w_\ell. \quad (9)$$

Following the proof of [Yan22, Thm. 1], (9) reduces to

$$\phi_i(x) = \sum_{\substack{\ell \in L \\ i \in \mathcal{F}_\ell}} \left( \sum_{S \subset \mathcal{F}_\ell \setminus \{i\}} \frac{|S|!(|\mathcal{F}_\ell| - 1 - |S|)!}{|\mathcal{F}_\ell|!} (W_{\ell,S \cup \{i\}} - W_{\ell,S}) \right) \cdot w_\ell. \quad (10)$$

Note that

$$\begin{aligned} W_{\ell,S \cup \{i\}} &= \prod_{\substack{n \in \mathcal{P}_\ell \\ j_n \in S \cup \{i\}}} o_n^{(\ell)} \prod_{\substack{n \in \mathcal{P}_\ell \\ j_n \notin S \cup \{i\}}} z_n^{(\ell)} \\ &= \prod_{\substack{n \in \mathcal{P}_\ell \\ j_n \in S}} o_n^{(\ell)} \prod_{\substack{n \in \mathcal{P}_\ell \\ j_n \notin S}} z_n^{(\ell)} \prod_{\substack{n \in \mathcal{P}_\ell \\ j_n = i}} o_n^{(\ell)} / z_n^{(\ell)} \\ &= W_{\ell,S} \cdot \prod_{\substack{n \in \mathcal{P}_\ell \\ j_n = i}} o_n^{(\ell)} / z_n^{(\ell)} \end{aligned}$$

and hence,

$$W_{\ell,S \cup \{i\}} - W_{\ell,S} = W_{\ell,S} \cdot t_i^{(\ell)},$$

which concludes the proof.  $\blacksquare$

For implementation aspects it is worth noting that  $W_{\ell,S}$  depends on  $S$  but not on  $i$ , while  $t_i^{(\ell)}$  depends on  $i$  but not on  $S$ .

*Remark 4.* The problem with the above version of the theorem in the privacy-preserving setting is that it is data dependent:

- i) The definition of  $o_n^{(\ell)}$  depends on the value of the predicate; Lemma 3 shows how to obliviously compute  $o_n^{(\ell)}$ .
- ii) The condition  $i \in \mathcal{F}_\ell$  depends on knowing  $\mathcal{F}_\ell$ , which we do not. We will therefore define a vector  $\mathbf{t}_\ell$  of size  $M$  with  $t_{\ell,i} = t_i^{(\ell)}$  if  $i \in \mathcal{F}_\ell$  and  $t_{\ell,i} = 0$  otherwise. Lemma 12 shows how to obliviously compute  $\mathbf{t}_\ell$ . We can then rewrite (8) as

$$\phi_i(x) = \sum_{\ell \in L} \left( \sum_{S \subset \mathcal{F}_\ell \setminus \{i\}} \frac{|S|!(|\mathcal{F}_\ell| - 1 - |S|)!}{|\mathcal{F}_\ell|!} \cdot W_{\ell,S} \right) \cdot t_{\ell,i} \cdot w_\ell. \quad (11)$$

Note that if a leaf  $\ell$  is such that  $i \notin \mathcal{F}_\ell$ , then the inner sum needs to exclude the summand where  $S = \mathcal{F}_\ell$  because  $(|\mathcal{F}_\ell| - 1 - |S|)!$  is not well-defined. This is addressed in Section 3.2.4.

- iii) To compute the inner sum, one needs to consider a (potentially) larger set  $\mathcal{R}_\ell \supset \mathcal{F}_\ell$  of size  $D$ , see Section 3.2.1, and sum over all  $2^D$  subsets of  $\mathcal{R}_\ell$ . We do not know the set  $\mathcal{R}_\ell$ , all we know is it contains  $\mathcal{F}_\ell$ . We will then make use of obliviously computed characteristic vectors to annihilate the contribution of subsets that violate the condition  $S \subset \mathcal{F}_\ell \setminus \{i\}$  (assuming  $i \in \mathcal{F}_\ell$  - otherwise  $t_{\ell,i}$  annihilates the inner sum altogether, see ii)).

The main contribution of this work is Theorem 2, a data-independent version of Theorem 1.

### 3.2.1 Size reduction.

SHAP values are originally defined as sums over subsets of  $\{1, \dots, M\}$  (see Definition 2) which introduces an overhead of  $2^M$  in the complexity of computing them. The main observation of **TreeSHAP**, Theorem 1, is that by changing the index of summation in (7) (i.e., summing over the leaves  $\ell \in L$  first), for each leaf  $\ell$ , it suffices to sum only over the subsets of the set  $\mathcal{F}_\ell$  (which has size at most  $D$ ). This essentially yields a plaintext version of the **TreeSHAP** algorithm. The challenge of designing a data-independent privacy-preserving algorithm is that we do not know the exact size of  $\mathcal{F}_\ell$ . Therefore, the best we can do in a privacy-preserving setting is to reduce the overhead to  $2^D$  for each leaf  $\ell$  (instead of  $2^{|\mathcal{F}_\ell|}$ ) by considering a possibly larger subset  $\mathcal{R}_\ell \supset \mathcal{F}_\ell$  of  $\{1, \dots, M\}$  (see below for the exact definitions) of cardinality exactly equal to  $D$  (i.e., a cardinality that is data independent). As  $D \ll M$  in most of the practical applications, this is still a major gain in complexity.

Our approach is based on sorting permutations and reduces the characteristic vector  $F_\ell$  of  $\mathcal{F}_\ell$  (a binary vector of size  $M$ ) to a vector  $R_\ell$  of size  $D$  that encodes the same information as  $F_\ell$ . Let  $\mathfrak{S}_M$  be the set of permutations in  $M$  letters and let  $\sigma_{F_\ell} \in \mathfrak{S}_M$  be a permutation that sorts (in ascending order) the binary vector  $F_\ell$ . Note that  $\sigma_{F_\ell}(F_\ell) \in \{0, 1\}^M$  consists of  $M - |\mathcal{F}_\ell|$  0's followed by  $|\mathcal{F}_\ell|$  1's. We denote by  $R_\ell \in \{0, 1\}^D$  the vector formed by the last  $D$  entries of  $\sigma_{F_\ell}(F_\ell)$ .

We call a *valid index* any of the  $|\mathcal{F}_\ell|$  indices in  $\mathcal{F}_\ell$  (equivalently, the position of the 1 in  $F_\ell$  respectively  $R_\ell$  it corresponds to). Otherwise, we call an index *invalid*; there are  $D - |\mathcal{F}_\ell|$  invalid indices in  $R_\ell$ , and we denote by  $\mathcal{I}_\ell = \mathcal{R}_\ell \setminus \mathcal{F}_\ell$  the set of invalid indices of  $\mathcal{R}_\ell$ . Note that the order of the valid indices in  $R_\ell$  depends on the choice of  $\sigma_{F_\ell}$  and hence, does not yield useful information about any natural order of the indices in the path (e.g., the top-to-bottom order of the nodes in the path).

The concept of valid indices will be used to encode the subsets  $S \subset \mathcal{F}_\ell$  in the inner sum of (8) in a way that does not require knowing  $\mathcal{F}_\ell$ . One way to get a data-independent algorithm for computing (8) is to rewrite the inner sum as a sum over all  $2^D$  subsets of  $\mathcal{R}_\ell$

and make sure the contribution to the sum of any subset containing invalid indices is 0. We formalize this intuition in Theorem 2 below after introducing the required notions.

**Lemma 6** (Complexity). *The  $2^{D-1}$  permutations  $\sigma_{\mathcal{F}_\ell} \in \mathfrak{S}_M$  are obtained by sorting the  $2^{D-1}$  binary vectors  $F_\ell$ . Note that  $\sigma_{\mathcal{F}_\ell}(F_\ell)$  is a by-product of the sorting algorithm from [BCD<sup>+</sup>23], and  $R_\ell$  can therefore be computed locally.*

### 3.2.2 Ordering the power set $P(\mathcal{R}_\ell)$ .

We now introduce a convenient linear ordering on the power set  $P(\mathcal{R}_\ell)$  of  $\mathcal{R}_\ell$  that is induced by the choice of  $\sigma_{\mathcal{F}_\ell}$  (i.e., is non-canonical) satisfying the following properties:

- i) The first subset in the ordering is the entire set  $\mathcal{R}_\ell$  and the last subset is the empty set.
- ii) The power set  $P(\mathcal{F}_\ell)$  corresponds to the last  $2^{|\mathcal{F}_\ell|}$  subsets whereas the first  $2^D - 2^{|\mathcal{F}_\ell|}$  subsets contain invalid indices.
- iii) The order of indexing of  $P(\mathcal{R}_\ell)$  depends on  $\sigma_{\mathcal{F}_\ell}$  and is therefore non-canonical; yet, ii) always holds.
- iv) The cardinality of a subset  $S \subset \mathcal{R}_\ell$  is recovered explicitly from the binary expansion of its index in the linear ordering.

Let  $\{\tau_1, \dots, \tau_D\}$  be the elements of  $\mathcal{R}_\ell$  with ordering determined by the permutation  $\sigma_{\mathcal{F}_\ell}$ . Note that the first  $D - |\mathcal{F}_\ell|$  indices are invalid and the last  $|\mathcal{F}_\ell|$  are valid. The subset corresponding to position  $s = 1, \dots, 2^D$  in the linear ordering of  $P(\mathcal{R}_\ell)$  is the one whose characteristic vector is equal to the binary expansion of  $2^D - s$  (msb-to-lsb binary representation). The cardinality of the subset can thus be obtained as the sum of the bits of the binary expansion. For instance, if  $D = 3$  then the ordering is the following:

$$\{\tau_1, \tau_2, \tau_3\}, \{\tau_1, \tau_2\}, \{\tau_1, \tau_3\}, \{\tau_1\}, \{\tau_2, \tau_3\}, \{\tau_2\}, \{\tau_3\}, \emptyset.$$

E.g., for  $s = 2$ , the binary expansion of  $8 - 2$  is 110 and hence the 2nd subset in the ordering is  $\{\tau_1, \tau_2\}$ . Note that this ordering is agnostic to the presence of invalid indices (which is secret information).

### 3.2.3 Computing $W_{\ell,S}$ .

Let  $O_\ell \in \{0, 1\}^D$  be the binary vector given by

$$O_{\ell,d} = \begin{cases} \prod_{\substack{n \in \mathcal{P}_\ell \\ j_n = \tau_d}} o_n^{(\ell)} & \text{if } \tau_d \text{ is a valid index,} \\ 1 & \text{otherwise.} \end{cases}$$

and let  $Z_\ell \in [0, 1]^D$  be the vector given by

$$Z_{\ell,d} = \begin{cases} \prod_{\substack{n \in \mathcal{P}_\ell \\ j_n = \tau_d}} z_n^{(\ell)} & \text{if } \tau_d \text{ is a valid index,} \\ 1 & \text{otherwise.} \end{cases}$$

**Lemma 7.** *We obviously compute  $O_\ell$  and  $Z_\ell$  as:*

$$O_\ell = \pi_D \circ \sigma_{\mathcal{F}_\ell} \left( \prod_{k=1}^D \left( o_{\ell,k-1} \cdot P_\ell^{(k)} + (1 - P_\ell^{(k)}) \right) \right)$$

and

$$Z_\ell = \pi_D \circ \sigma_{F_\ell} \left( \prod_{k=1}^D \left( z_{\ell, k-1} \cdot P_\ell^{(k)} + (1 - P_\ell^{(k)}) \right) \right),$$

respectively, where  $P_\ell^{(k)}$  is the  $k$ th column of  $P_\ell$  and  $\pi_D: \mathbb{R}^M \rightarrow \mathbb{R}^D$  is the projection onto the last  $D$  coordinates.

*Proof.* Let  $\tau_d$  be a valid index, that is,  $T_d := \{1 \leq k \leq D: j_{n, k-1}^{(\ell)} = \tau_d\}$  is non-empty.

Moreover, for all  $k \in T_d$  we have  $P_\ell^{(k)} = \mathbf{b}_{\tau_d}$  (the only non-zero entry is at position  $\tau_d$ ), and for all  $k \in \{1, \dots, D\} \setminus T_d$  the  $\tau_d$ th entry of  $P_\ell^{(k)}$ , which we denote by  $P_{\ell, \tau_d}^{(k)}$ , is 0. Hence,

$$\begin{aligned} \prod_{\substack{n \in \mathcal{P}_\ell \\ j_n = \tau_d}} o_n^{(\ell)} &= \prod_{k \in T_d} o_{\ell, k-1} \\ &= \prod_{k=1}^D \left( o_{\ell, k-1} \cdot P_{\ell, \tau_d}^{(k)} + (1 - P_{\ell, \tau_d}^{(k)}) \right), \end{aligned}$$

and  $\pi_D \circ \sigma_{F_\ell}$  sends the  $\tau_d$ th entry of a vector of size  $M$  to the  $d$ th entry of a vector of size  $D$ .

For an invalid index  $\tau_d$  one observes that  $P_{\ell, \tau_d}^{(k)} = 0$  for all  $k = 1, \dots, D$  and hence,

$$\prod_{k=1}^D \left( o_{\ell, k-1} \cdot P_{\ell, \tau_d}^{(k)} + (1 - P_{\ell, \tau_d}^{(k)}) \right) = 1,$$

which concludes the proof for  $O_\ell$ . The proof for  $Z_\ell$  is identical.  $\blacksquare$

*Remark 5.* The (plaintext) definition (5) of  $W_{\ell, S}$  applies to subsets  $S \subset \mathcal{F}_\ell$ . We can naturally extend this definition to  $S \subset \mathcal{R}_\ell$  by setting  $W_{\ell, S} := W_{\ell, S \cap \mathcal{F}_\ell}$  and note that Lemma 5 still holds.

Let  $W_\ell^{P(\mathcal{R}_\ell)} \in [0, 1]^{2^D}$  be the vector whose  $S$ th entry is  $W_{\ell, S}$  (when indexing the power set  $P(\mathcal{R}_\ell)$  by the subsets of  $\mathcal{R}_\ell$  as described in Section 3.2.2).

**Lemma 8.** *We obviously compute  $W_\ell^{P(\mathcal{R}_\ell)}$  as:*

$$W_\ell^{P(\mathcal{R}_\ell)} = \prod_{d=1}^D \text{concat}_d(O_{\ell, d}, Z_{\ell, d}),$$

where  $\text{concat}_d(a, b) = \underbrace{(a, \dots, a)}_{2^{D-d}} \underbrace{(b, \dots, b)}_{2^{D-d}} \dots \underbrace{(a, \dots, a)}_{2^{D-d}} \underbrace{(b, \dots, b)}_{2^{D-d}} \in \mathbb{R}^{2^D}$ .

*Proof.* According to Lemma 5 and Remark 5, for any subset  $S \subset \mathcal{R}_\ell$ ,

$$W_{\ell, S} = \prod_{\substack{n \in \mathcal{P}_\ell \\ j_n \in S}} o_n^{(\ell)} \prod_{\substack{n \in \mathcal{P}_\ell \\ j_n \notin S}} z_n^{(\ell)}.$$

Fix a subset  $S \subset \mathcal{R}_\ell$  with corresponding position  $1 \leq s \leq 2^D$ , and let  $b_{D-1} \dots b_0$  be the binary expansion of  $2^D - s$ . Recall that  $\tau_1 \in S \Leftrightarrow b_{D-1} = 1$ ,  $\tau_2 \in S \Leftrightarrow b_{D-2} = 1$ , etc. Since  $O_{\ell, d} = \prod_{\substack{n \in \mathcal{P}_\ell \\ j_n = \tau_d}} o_n^{(\ell)}$  and  $Z_{\ell, d} = \prod_{\substack{n \in \mathcal{P}_\ell \\ j_n = \tau_d}} z_n^{(\ell)}$  if  $\tau_d$  is a valid index (and  $O_{\ell, d} = Z_{\ell, d} = 1$  if  $\tau_d$  is an invalid index), one deduces that

$$W_{\ell, S} = \prod_{d=1}^D (b_{D-d} ? O_{\ell, d} : Z_{\ell, d}).$$

Moreover, the  $s$ th entry of  $\text{concat}_d(a, b)$  equals  $a$  if and only if  $(2^D - s) / 2^{D-d}$  is odd (Euclidean division) if and only if  $b_{D-d} = 1$ , which concludes the proof.  $\blacksquare$

**Lemma 9** (Complexity).

- There are  $2^D$  binary vectors  $O_\ell$  of size  $D \times 1$ , obtained via Lemma 7, requiring  $D$  ANDs between a bit and a binary vector of size  $M \times 1$ ,  $D - 1$  ANDs between binary vectors of size  $M \times 1$  and one call to  $\sigma_{F_\ell}$  each.
- There are  $2^D$  vectors  $Z_\ell$  of size  $D \times 1$ , obtained via Lemma 7, requiring  $D$  multiplications between a scalar and a binary vector of size  $M \times 1$ ,  $D - 1$  multiplications between vectors of size  $M \times 1$  and one call to  $\sigma_{F_\ell}$  each.
- There are  $2^D$  vectors  $W_\ell^{P(\mathcal{R}_\ell)}$  of size  $2^D \times 1$ , obtained via Lemma 8, requiring  $D - 1$  multiplications between vectors of size  $2^D \times 1$  each.

### 3.2.4 Computing Shapley weights.

To filter subsets that contain valid indices only, we need to make the following definition:

**Definition 2.** Let  $S \subset \mathcal{R}_\ell$ . The *Shapley weight* of  $S$  with respect to  $\mathcal{F}_\ell$  is

$$C_S = \begin{cases} \frac{|S|!(|\mathcal{F}_\ell| - 1 - |S|)!}{|\mathcal{F}_\ell|!} & \text{if } S \subsetneq \mathcal{F}_\ell \\ 0 & \text{otherwise.} \end{cases}$$

*Remark 6.* The reason why we exclude  $S = \mathcal{F}_\ell$  is the inner summation of (11). If  $i \in \mathcal{F}_\ell$  then  $S \subset \mathcal{F}_\ell \setminus \{i\}$  implies  $S \subsetneq \mathcal{F}_\ell$ , and if  $i \notin \mathcal{F}_\ell$  then the case  $S = \mathcal{F}_\ell$  would yield factorial of a negative number.

Since  $\mathcal{F}_\ell$  and  $|\mathcal{F}_\ell|$  are not known a priori, we first compute a larger public matrix  $C$  of size  $2^D \times D$  whose  $d$ th column  $C^{(d)}$  is the vector of the  $2^D$  Shapley weights (one for each subset  $S \subset \mathcal{R}_\ell$ ) for the case  $|\mathcal{F}_\ell| = d$ . In a second step we compute  $\delta_{|\mathcal{F}_\ell|}$ , the secret characteristic vector of  $|\mathcal{F}_\ell|$ , and obviously select the correct column of  $C$ .

We compute the  $s$ th entry  $C_s^{(d)}$  of the  $d$ th column of  $C$  as

$$C_s^{(d)} = \begin{cases} \frac{|S|!(d-1-|S|)!}{d!} & \text{if } s > 2^D - 2^d + 1 \\ 0 & \text{otherwise,} \end{cases}$$

for all  $d = 1, \dots, D$  and all  $s = 1, \dots, 2^D$ , where the cardinality of the corresponding subset  $S$  at position  $s$  is given by property iv) from Section 3.2.2. Note that the matrix  $C$  is pre-computed (it is the same for all  $i \in \{1, \dots, M\}$  and for all  $\ell \in L$ ).

We are now interested in obviously computing the secret-shared, one-hot encoding (indicator) vector  $\delta_{|\mathcal{F}_\ell|}$  of the secret cardinality  $|\mathcal{F}_\ell|$ , that is, the binary vector of size  $D$  that has a single 1 at the  $|\mathcal{F}_\ell|$ th position. This is done via the following simple lemma:

**Lemma 10.** If  $R_{\ell,d}$  denotes the  $d$ th entry of the vector  $R_\ell = (\underbrace{0, \dots, 0}_{D-|\mathcal{F}_\ell|}, \underbrace{1, \dots, 1}_{|\mathcal{F}_\ell|})$  then

$$\delta_{|\mathcal{F}_\ell|} = (R_{\ell,D} \oplus R_{\ell,D-1}, \dots, R_{\ell,2} \oplus R_{\ell,1}, R_{\ell,1})^t \in \{0, 1\}^D.$$

Since we have  $R_\ell$  available in secret shares, the above lemma yields secret shares for  $\delta_{|\mathcal{F}_\ell|}$ . Finally, we compute  $C^{(|\mathcal{F}_\ell|)}$ , the  $|\mathcal{F}_\ell|$ th column of  $C$ , as

$$C^{(|\mathcal{F}_\ell|)} = \sum_{d=1}^D \delta_{|\mathcal{F}_\ell|,d} \cdot C^{(d)} \in [0, 1]^{2^D}. \quad (12)$$

**Lemma 11** (Complexity).

- $C$  is a publicly computed matrix of size  $2^D \times D$ .
- There are  $2^D$  binary vectors  $\delta_{|\mathcal{F}_\ell|}$  of size  $D \times 1$ , obtained via Lemma 10, which can be computed locally from the binary vector  $R_\ell$ .
- There are  $2^D$  vectors  $C^{(|\mathcal{F}_\ell|)}$  of size  $2^D \times 1$ , which are computed locally as the matrix multiplication between the public matrix  $C$  and the secret vector  $\delta_{|\mathcal{F}_\ell|}$ .

### 3.2.5 Computing $\mathbf{t}_\ell$ .

Let  $\mathbf{t}_\ell \in [-1, N - 1]^M$  be the vector with  $i$ th entry given by

$$t_{\ell,i} = \begin{cases} \prod_{\substack{n \in \mathcal{P}_\ell \\ j_n = i}} o_n^{(\ell)} / z_n^{(\ell)} - 1 & \text{if } i \in \mathcal{F}_\ell \\ 0 & \text{otherwise,} \end{cases}$$

where  $N$  is the number of samples the tree model was trained with.

**Lemma 12.** *We obviously compute  $\mathbf{t}_\ell$  as:*

$$\mathbf{t}_\ell = \prod_{k=1}^D \left( o_{\ell,k-1} \cdot z_{\ell,k-1}^{-1} \cdot P_\ell^{(k)} + (1 - P_\ell^{(k)}) \right) - 1,$$

where  $P_\ell^{(k)}$  is the  $k$ th column of  $P_\ell$ .

*Proof.* For all  $i = 1, \dots, M$ , let  $T_i := \{1 \leq k \leq D : j_{n_{k-1}}^{(\ell)} = i\}$ . For all  $k \in T_i$  we have  $P_\ell^{(k)} = \mathbf{b}_i$  (the only non-zero entry is at position  $i$ ), and for all  $k \in \{1, \dots, D\} \setminus T_i$  the  $i$ th entry of  $P_\ell^{(k)}$ , which we denote by  $P_{\ell,i}^{(k)}$ , is 0. Hence, if  $i \in \mathcal{F}_\ell$  then

$$\begin{aligned} \prod_{\substack{n \in \mathcal{P}_\ell \\ j_n = i}} o_n^{(\ell)} / z_n^{(\ell)} - 1 &= \prod_{k \in T_i} o_{\ell,k-1} \cdot z_{\ell,k-1}^{-1} - 1 \\ &= \prod_{k=1}^D \left( o_{\ell,k-1} \cdot z_{\ell,k-1}^{-1} \cdot P_{\ell,i}^{(k)} + (1 - P_{\ell,i}^{(k)}) \right) - 1. \end{aligned}$$

And if  $i \notin \mathcal{F}_\ell$  then  $T_i = \emptyset$  and hence,

$$\prod_{k=1}^D \left( o_{\ell,k-1} \cdot z_{\ell,k-1}^{-1} \cdot P_{\ell,i}^{(k)} + (1 - P_{\ell,i}^{(k)}) \right) - 1 = \prod_{k=1}^D 1 - 1 = 0. \quad \blacksquare$$

**Lemma 13** (Complexity).

- Computing all inverse cover ratios  $1/z_n^{(\ell)}$  in Lemma 12 requires two private divisions between vectors of size  $2^D - 1 \times 1$  (unlike for the computation of  $\mathbf{z}_\ell$ , there is no easily computable relation between the inverse cover ratio of a left child node and the inverse cover ratio of its right sibling node).
- There are  $2^D$  vectors  $\mathbf{t}_\ell$  of size  $M \times 1$ , obtained via Lemma 12, requiring  $D$  multiplications between two scalars,  $D$  multiplications between a scalar and a binary vector of size  $M \times 1$  and  $D - 1$  multiplications between vectors of size  $M \times 1$  each.

### 3.2.6 Putting everything together.

For  $i \in \{1, \dots, M\}$ , let  $\mathbf{I}_\ell^{(i)} \in \{0, 1\}^{2^D}$  be the indicator vector of  $\{S \subset \mathcal{R}_\ell : i \notin S\}$ , i.e.,  $\mathbf{I}_{\ell,s}^{(i)} = 1$  if and only if  $i \notin S$ , where  $S \subset \mathcal{R}_\ell$  is the subset corresponding to position  $s$ .

**Lemma 14.** *We have*

$$\mathbf{I}_\ell^{(i)} = \neg \left( \bigoplus_d^D \text{concat}_d(\eta_{\ell,d}^{(i)}, 0) \right),$$

$$\text{where } \eta_\ell^{(i)} = \pi_D \circ \sigma_{F_\ell} \left( \underbrace{(0, \dots, 0)}_{i-1}, F_{\ell,i}, \underbrace{(0, \dots, 0)}_{M-i} \right) \in \{0, 1\}^D.$$

**Theorem 2.** *We obviously compute  $\phi_i(x)$  as follows:*

$$\phi_i(x) = \sum_{\ell \in L} \left( \sum_{s=1}^{2^D} \mathbf{I}_{\ell,s}^{(i)} \cdot C_s^{(|\mathcal{F}_\ell|)} \cdot W_{\ell,s}^{P(\mathcal{R}_\ell)} \right) \cdot t_{\ell,i} \cdot w_\ell. \quad (13)$$

*Proof.* We want to show that (13) computes (8). For all leaves  $\ell \in L$  we have  $t_{\ell,i} = 0$  whenever  $i \in \{1, \dots, M\} \setminus \mathcal{F}_\ell$  and hence, (8) reduces to

$$\sum_{\ell \in L} \left( \sum_{S \subset \mathcal{F}_\ell \setminus \{i\}} \frac{|S|!(|\mathcal{F}_\ell| - 1 - |S|)!}{|\mathcal{F}_\ell|!} \cdot W_{\ell,S} \right) \cdot t_{\ell,i} \cdot w_\ell.$$

Regarding the inner sum, for each leaf  $\ell \in L$ , the subsets  $S \subsetneq \mathcal{F}_\ell$  correspond to the positions  $s = 2^D - 2^{|\mathcal{F}_\ell|} + 2, \dots, 2^D$  with respect to the ordering defined in Section 3.2.2. Moreover,  $C_s^{(|\mathcal{F}_\ell|)} = 0$  for  $s = 1, \dots, 2^D - 2^{|\mathcal{F}_\ell|} + 1$  and  $C_s^{(|\mathcal{F}_\ell|)} = \frac{|S|!(|\mathcal{F}_\ell| - 1 - |S|)!}{|\mathcal{F}_\ell|!}$  for  $s = 2^D - 2^{|\mathcal{F}_\ell|} + 2, \dots, 2^D$ . Note that the Shapley weight depends on the cardinality of  $S$  only, which by iv) of Section 3.2.2 is uniquely determined by the position  $s$  in the ordering (and not the ordering itself). Since  $W_\ell^{P(\mathcal{R}_\ell)}$  is ordered in a way that respects iv) of Section 3.2.2, for each subset  $S \subsetneq \mathcal{F}_\ell$  the quantity  $W_{\ell,S}$  is multiplied with the right Shapley weight  $\frac{|S|!(|\mathcal{F}_\ell| - 1 - |S|)!}{|\mathcal{F}_\ell|!}$ .

Finally, the characteristic vector  $\mathbf{I}_\ell^{(i)}$  annihilates the contribution of all subsets  $S$  that contain  $i$ , which concludes the proof.  $\blacksquare$

**Lemma 15** (Complexity).

- Computing the  $M2^{D-1}$  binary vectors  $\eta_\ell^{(i)}$  and  $\mathbf{I}_\ell^{(i)}$  respectively from Lemma 14 requires  $M2^{D-1}$  permutations of binary vectors of size  $M \times 1$ , all other computations can be performed locally.
- Computing  $\phi_i(x)$  from Theorem 2: the inner sum requires  $2^D \cdot 2$  scalar multiplications, hence  $\phi_i(x)$  is computed with a total of  $2^D(2^{D+1} + 2)$  scalar multiplications.

### 3.3 XorSHAP for the 0-cover case.

As pointed out in Remark 1, the definition of  $W_{\ell,S}(x)$ , see (5), is only valid if all split nodes  $n$  have positive cover (i.e.,  $c_n > 0$ ). In a plaintext training algorithm, any split node with 0-cover will be removed during pruning. This, however, modifies the tree structure. A full binary tree of depth  $D$  is a tree with  $2^D - 1$  split nodes and  $2^D$  leaf nodes. In order for XorSHAP to work correctly in the presence of 0-cover nodes, the full binary trees of the model are supposed to be pruned during training in a way that preserves the tree structure.

By *pruning a node* of a decision tree we mean the process of replacing the subtree rooted at that node by a leaf node with some specified weight.

**Definition 3.** Let **Pruning** be any plaintext pruning algorithm. Its tree structure-preserving counterpart **PruningFBT** is the following algorithm: suppose that **Pruning** replaces the subtree rooted at a node  $n$  by a single leaf node with some weight  $w$ ; then, **PruningFBT** is the algorithm that keeps the binary subtree rooted at  $n$  and only replaces all leaf weights of that subtree by  $w$ .

A pruned full binary tree by a **PruningFBT** is thus still a full binary tree, however, as soon as a sample evaluates through a split node  $n$  that got pruned the predicted value no longer depends on the taken sub-path starting from  $n$ . While evaluating a FBT-pruned full binary tree **Tree** that contains 0-cover nodes is straightforward, computing the conditional expectation  $\mathbb{E}[\mathbf{Tree}(\xi) \mid \xi_S = x_S]$  is not.

Let **Tree** be a pruned full binary tree (it may contain 0-cover nodes). We define  $G'_S(x; n)$ , a generalization of  $G_S(x; n)$ , see (4), as

$$G'_S(x; n) = \begin{cases} w_n & \text{if } n \text{ is a leaf node,} \\ G'_S(x; n_{\text{right}}) & \text{else if } c_{n_{\text{left}}} = 0, \\ G'_S(x; n_{\text{left}}) & \text{else if } c_{n_{\text{right}}} = 0, \\ (x_{j_n} < t_n) ? G'_S(x; n_{\text{left}}) : G'_S(x; n_{\text{right}}) & \text{else if } j_n \in S, \\ \frac{c_{n_{\text{left}}}}{c_n} G'_S(x; n_{\text{left}}) + \frac{c_{n_{\text{right}}}}{c_n} G'_S(x; n_{\text{right}}) & \text{else.} \end{cases} \quad (14)$$

**Lemma 16.** Let **Tree** be a full binary tree containing 0-cover nodes and let **Pruning** be a plaintext pruning algorithm, with **PruningFBT** its tree structure preserving counterpart. If  $G'_S(x; n_{\text{root}})$  is computed with **PruningFBT(Tree)** then

$$\mathbb{E}[\mathbf{Pruning}(\mathbf{Tree})(\xi) \mid \xi_S = x_S] = G'_S(x; n_{\text{root}}).$$

*Proof.* Suppose that sample  $x$  evaluates through a split node  $n$  with non-zero cover, however a child node of  $n$  has zero cover. Without loss of generality assume that  $c_{n_{\text{left}}} = 0$ . Since **Pruning** replaces  $n$  by a leaf node (say with weight  $w$ ), all leaf nodes of the full subtree of **PruningFBT(Tree)** with root  $n$  will have weight  $w$ . A simple induction then shows that  $G'_S(x; n_{\text{right}}) = w$ . ■

We can extend Definition 1 to a notion of consistency between a path  $\mathcal{P}_\ell$  and a pair  $(x, S)$  of a sample and a feature subset to the case where a full binary tree contains 0-cover nodes.

**Definition 4.** We call a path  $\mathcal{P}_\ell$  *consistent* with respect to the pair  $(x, S)$  if every node in  $\mathcal{P}_\ell \cup \{\ell\}$  has non-zero cover and if for every node  $n \in \mathcal{P}_\ell$  such that  $j_n \in S$ , either the child node determined by the condition  $x_{j_n} < t_n$  coincides with the node of  $\mathcal{P}_\ell \cup \{\ell\}$  that succeeds  $n$  in the path or the child node determined by the condition  $\neg(x_{j_n} < t_n)$  has zero cover. Otherwise, the path is called *inconsistent*.

The quantity  $W_{\ell, S}$ , see (5), can naturally be extended to the case of full binary trees containing 0-cover nodes, with the notion of consistency introduced in Definition 4.

It remains to show how to obviously evaluate  $G'_S(x; n)$  for a pruned oblivious (i.e., full binary) tree **Tree**. Note that only the weights of leaves whose path is consistent with  $(x, S)$  contribute to  $G'_S(x; n)$ . Unlike  $\mathbf{eval}(x, \mathbf{Tree})$ , which can be obviously computed from the secret bits  $\beta = (x_j < t)$ , see [DDG<sup>+</sup>22], one needs to consider more than the bits  $\beta$  to evaluate  $G'_S(x; n)$ . Recall that for the evaluation of  $G_S(x; n)$ , one computes  $\mathbf{o}_\ell$  from the bits  $\beta$  only, see Lemma 3.

Given a leaf node  $\ell \in L$ , let  $\mathcal{P}_\ell$  be the path of split nodes leading to  $\ell$ . We define  $\tilde{\mathbf{o}}_\ell = (\tilde{o}_n^{(\ell)} : n \in \mathcal{P}_\ell)$ , similar to  $\mathbf{o}_\ell$ , where  $\tilde{o}_n^{(\ell)} = 1$  if the successor node  $n_{\text{child}}$  of  $n$  in  $\mathcal{P}_\ell \cup \{\ell\}$

has non-zero cover and either evaluating node  $n$  on the sample  $x$  yields  $n_{\text{child}}$  or the sibling node of  $n_{\text{child}}$  (when looking at the whole tree) has zero cover; otherwise  $\tilde{o}_n^{(\ell)} = 0$ .

**Lemma 17.**

$$W_{\ell,S}(x) = \prod_{\substack{n \in \mathcal{P}_\ell \\ j_n \in S}} \tilde{o}_n^{(\ell)} \prod_{\substack{n \in \mathcal{P}_\ell \\ j_n \notin S}} z_n^{(\ell)}. \quad (15)$$

*Proof.* If  $\mathcal{P}_\ell \cup \{\ell\}$  contains a 0-cover node then the path is inconsistent with  $(x, S)$ . In this case, let  $n$  be the unique node in the path with non-zero cover and whose successor node  $n_{\text{child}}$  has zero cover (a child of a 0-cover node has zero cover, so  $n$  is unique in the path). If  $j_n \in S$  then  $\tilde{o}_n^{(\ell)} = 0$ , else  $z_n^{(\ell)} = \frac{c_{n_{\text{child}}}}{c_n} = 0$ .

If  $\mathcal{P}_\ell \cup \{\ell\}$  does not contain any 0-cover node then  $\mathcal{P}_\ell$  is consistent with  $(x, S)$  if and only if  $\prod_{\substack{n \in \mathcal{P}_\ell \\ j_n \in S}} \tilde{o}_n^{(\ell)} = 1$ . ■

Note that for any 0-cover split node  $n$  the ratio  $z_n^{(\ell)}$  is not well-defined. A data-oblivious division algorithm is expected to return arbitrary output if the denominator is 0 and hence, the presence of 0-cover split nodes in the path does not invalidate Lemma 17.

It remains to show how to obviously compute  $\tilde{o}_\ell$ . One then applies Lemma 7 and Lemma 8 with  $\tilde{o}_\ell$  to obviously compute  $\phi_i(x)$  as in Theorem 2. Fix an index  $d$ . As in Lemma 3, let  $(j, t)$  be the feature index and threshold associated with node  $n_d^{(\ell)}$ , and let  $\beta = (x_j < t) \in \{0, 1\}$ . Let  $n_{\text{left}}$  and  $n_{\text{right}}$  be the left and right children nodes of  $n_d^{(\ell)}$  respectively (when looking at the whole tree). Consider the following bits

$$\gamma_{\text{left}} = ((c_{n_{\text{left}}} == 0) ? 0 : ((c_{n_{\text{right}}} == 0) ? 1 : \beta))$$

and

$$\gamma_{\text{right}} = \neg\gamma_{\text{left}} \oplus (c_{n_{\text{left}}} == 0) \wedge (c_{n_{\text{right}}} == 0),$$

which can be computed obliviously.

**Lemma 18.**

$$\tilde{o}_{\ell,d} = \begin{cases} \gamma_{\text{left}} & \text{if } n_{\text{child}} = n_{\text{left}} \\ \gamma_{\text{right}} & \text{if } n_{\text{child}} = n_{\text{right}} \end{cases}$$

*Proof.* We have:

- if  $c_{n_{\text{left}}} > 0$  and  $c_{n_{\text{right}}} > 0$  then  $\gamma_{\text{left}} = \beta$  and  $\gamma_{\text{right}} = \neg\beta$ ,
  - if  $c_{n_{\text{left}}} > 0$  and  $c_{n_{\text{right}}} = 0$  then  $\gamma_{\text{left}} = 1$  and  $\gamma_{\text{right}} = 0$ ,
  - if  $c_{n_{\text{left}}} = 0$  and  $c_{n_{\text{right}}} > 0$  then  $\gamma_{\text{left}} = 0$  and  $\gamma_{\text{right}} = 1$ ,
  - if  $c_{n_{\text{left}}} = 0$  and  $c_{n_{\text{right}}} = 0$  then  $\gamma_{\text{left}} = 0$  and  $\gamma_{\text{right}} = 0$ .
- 

## 4 Complexity Analysis

For analyzing the complexity of XorSHAP, consider a model consisting of a single tree. For an arbitrary tree ensemble model, one can simply multiply the complexity by the number  $T$  of decision trees in the ensemble.

Below, we list the complexities of the computations needed in the different sections of Section 3.

- §3.1 By Lemma 4, the cost is dominated by the computation of the  $\mathbf{z}_\ell$ 's, which amounts to one private division between two vectors of size  $2^D - 1 \times 1$ .
- §3.2.1 By Lemma 6, the cost is dominated by the computation of the  $\sigma_{\mathcal{F}_\ell}$ 's, which amounts to the sorting of  $2^{D-1}$  binary vectors of size  $M$ .
- §3.2.3 By Lemma 9, the cost is dominated by the computation of the  $Z_\ell$ 's, which amounts to  $2^D$  times  $D-1$  multiplications between vectors of size  $M \times 1$ , or by the computation of the  $W_\ell^{P(\mathcal{R}_\ell)}$ 's, which amounts to  $2^D$  times  $D-1$  multiplications between vectors of size  $2^D \times 1$ .
- §3.2.4 By Lemma 11, all computations happen locally.
- §3.2.5 By Lemma 13, the cost is dominated by the computation of the  $\mathbf{t}_\ell$ 's, which amounts to  $2^D$  times  $D-1$  multiplications between vectors of size  $M \times 1$ .
- §3.2.6 By Lemma 15, the cost is dominated by the computation of  $\phi_i(x)$ , which amounts to  $2^D(2^{D+1} + 2)$  scalar multiplications.

Comparison of the above-listed complexities shows that the bottlenecks are:

- Computing all  $Z_\ell$ 's from Section 3.2.3 is  $O(MD2^D)$ ,
- Computing all  $W_\ell^{P(\mathcal{R}_\ell)}$ 's from Section 3.2.3 is  $O(D2^{2D})$ ,
- Computing all  $\mathbf{t}_\ell$ 's from Section 3.2.5 is  $O(MD2^D)$ ,
- Computing  $\phi_i(x)$  from Section 3.2.6 is  $O(2^{2D+1})$ .

Hence, if  $\tilde{M} = \max(M, 2^D)$ , then **XorSHAP** runs in

$$O(T\tilde{M}D2^D)$$

multiplications, where  $T$  is the number of trees.

In contrast to **TreesSHAP** (see Section 2.2.3), whose complexity is  $O(T2^D D^2)$  (i.e., only depends on the depth  $D$  and not on the number of features  $M$ ), **XorSHAP**'s runtime is linear in  $M$ . This is the price we pay for having a data oblivious algorithm since feature subsets are represented by characteristic vectors of size  $M$  (though we avoid having overhead  $2^M$  by reducing the representation of feature subsets to characteristic vectors of size  $D$ , see Section 3.2.1).

## 5 Implementation and Benchmarks

To demonstrate the scalability of **XorSHAP** we implemented the algorithm on the **Manticore** framework [BCD<sup>+</sup>23] and ran benchmarks for various input dimensions and model parameters. In our experiments, we are only interested in the runtime/resources (though we check correctness vs standard implementations). For this purpose, the trees are randomly generated full binary trees of depth  $D$ , that is, for each split node of each tree we generate a random feature index  $j \in \{1, \dots, M\}$  and a random threshold  $t \in \mathbb{R}$ , and for each leaf node we generate a random weight  $w \in \mathbb{R}$  and a random cover  $c \geq 0$  (the sum of the covers of all leaf nodes was fixed to 10'000). The test dataset is a randomly generated dataset of desired dimensions.

We reiterate that the sole purpose of the presented benchmarks (run with the random input models and data) was to provide practical complexity estimates. As described in Section 6, we have already used **XorSHAP** in real-world settings, models and input data and with meaningful implications of privacy-preserving model explainability.

The offline phase (generation of random precomputed data) is performed by the trusted dealer. The benchmarks were run in a 2-party scenario with input data already provided in secret shares across the two parties (i.e., simulating any kind of partition). We vary the number of trees  $T$ , the tree depth  $D$ , the number of features  $M$  as well as the number of testing samples (simultaneous computation of the SHAP values). The benchmarks include offline and online phase and were run on an Intel Xeon E5-2666 v3 @ 2.90GHz CPU with 32GB RAM. Over all runs, the  $L_\infty$ -distance between the SHAP values computed in MPC and the baseline plaintext computation with the classical SHAP implementation<sup>2</sup> is around  $10^{-13}$ .

In Table 1 and Table 2, we provide benchmarks (offline/online runtime and number of outgoing connections) and (network transfer per party and memory overhead) respectively, across different values of  $T$  for simultaneous computations of 20 samples/predictions, for  $M = 20$  features and binary trees of depth  $D = 6$ . Similarly, in Table 3 and Table 4, we provide the same type of data for varying number of features  $M$ .

Figure 2 shows the total runtime for varying number of trees and tree depth, and varying number of features and simultaneously evaluated samples respectively. The time on the  $y$ -axis consists of the offline compute time plus the online compute time plus and estimate of the network transfer time. The latter is computed by using the network transfer (offline triplet size and online network size), number of outgoing connections, the throughput and the latency of the network. The estimation uses a network whose throughput is 120Mbps and whose latency is 0.3ms.

**Table 1:** #samples = 20,  $M = 20$ ,  $D = 6$

T	offline compute time (in s)	online compute time (in s)	number of outgoing connections
100	97	95	334
80	75	74	334
60	57	55	334
40	37	36	334
20	17	14	334

**Table 2:** #samples = 20,  $M = 20$ ,  $D = 6$

T	network transfer (in MB, per party, offline + online)	memory offline (in MB)	memory online (in MB)
100	5329	3867	11679
80	4263	3112	9379
60	3196	2346	7056
40	2131	1572	4744
20	1065	806	2430

**Table 3:**  $D = 4$ ,  $T = 60$ , #samples = 100

M	offline compute time (in s)	online compute time (in s)	number of outgoing connections
100	64	60	399
80	50	49	379
60	38	38	356
40	26	25	336
20	15	14	313

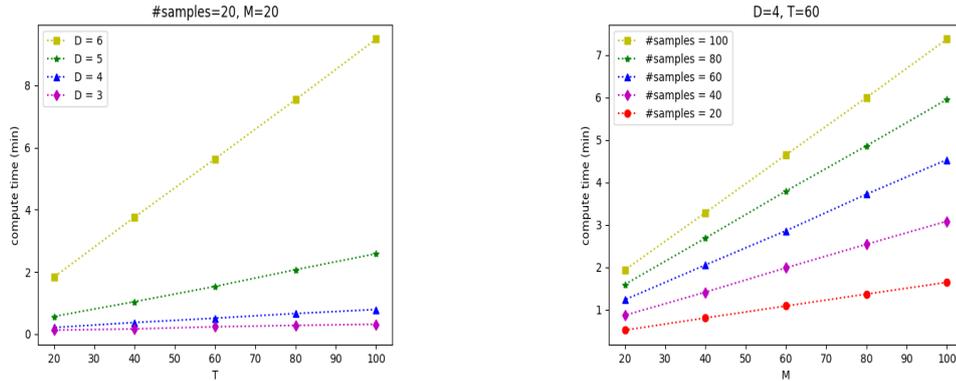
## 6 Privacy Implications in real-world use cases

While our benchmarks presented in the previous section were computed on random input data, our algorithm has already been applied to many real-world use cases in financial

<sup>2</sup>See <https://shap.readthedocs.io/en/latest/> as well as the function `shap.TreeExplainer.shap_values(x)`.

**Table 4:**  $D = 4$ ,  $T = 60$ , #samples = 100

M	network transfer (in MB, per party, offline + online)	memory offline (in MB)	memory online (in MB)
100	4606	2834	5836
80	3751	2277	4691
60	2897	1739	3554
40	2043	1191	2444
20	1189	632	1295

**Figure 2:** Left: varying number of trees and tree depth. Right: varying number of features and simultaneously evaluated samples.

technologies, healthcare as well as predictive maintenance where the input data and the model are sensitive and need to be protected.

We now summarize one important application in root-cause analysis and yield optimization in the semiconductor fabrication process of microchips - a process that includes two major parties, the Fabs (fabrication plants that manufacture semiconductor products either for other companies or for themselves) and the SMEs (Semiconductor Manufacturing Equipment companies that produce advanced and highly specialized equipment for the manufacturing process). The input data consists of sensor readings from the manufacturing equipment (accessible to the SMEs) and labels (typically assigned by the Fabs). One uses supervised machine learning models, very often tree ensemble models, to predict the labels from the sensor readings data. Model explainability is crucial to identify the input features (sensor readings) that affects a predicted defect of a wafer.

From a privacy perspective, the Fabs and the SMEs are not willing to share their data with each other due to sensitivity - Fabs are interested to protect their manufacturing pipelines whereas SMEs are reluctant to give their sensor readings to the Fabs due to IP concerns.

As there are multiple sensor readings in the prediction model, a lot of these are irrelevant for explaining the predicted defect. Note that **XorSHAP** enables privacy-preserving computation of SHAP values (as explainability metrics) to better identify the relevant sensor readings and therefore, the root cause of the defect. We refer the interested reader to [JSTV23] for more details on this and other use cases.

Note that the above real-world use cases are already within reach of the current implementation where **XorSHAP** already supports decision tree ensemble models with hundreds of trees (of reasonable depth  $D \geq 10$ ) as well as hundreds of features  $M$  (reflecting the complexity of the manufacturing process in terms of sensor readings).

## 7 Conclusion

The proposed XorSHAP is one of the first attempt to compute SHAP values in an efficient and scalable fashion for decision tree ensemble models in a generic data-oblivious manner, thus, making it applicable to various privacy-preserving settings such as SMPC, FHE as well as differential privacy (DP) as well as to a multitude of choices on the underlying threat models and security assumptions. On a decision tree ensemble model with  $T$  binary decision trees, with  $M$  features in depth  $D$ , our algorithm has a run-time  $O(T\widetilde{M}D^D)$  where  $\widetilde{M} = \max(M, 2^D)$ , thus, being linear in the number of features.

We implemented XorSHAP using the Manticore framework for Secure Multiparty Computation (SMPC) with additive secret sharing using modular real secret shares and Boolean secret shares. The security model of our implementation, semi-honest with full-threshold and trusted dealer in the offline phase is well-adapted to a lot of the practical use cases. The implementation scales to real-world datasets (e.g., computing SHAP values for models with hundreds of trees, each with hundreds of features) and is already part of Inpher XOR Platform for privacy-preserving computations addressing relevant real-world use cases as described in Section 6.

Finally, the parallelization-friendly nature of our method opens another future research direction on massively scaling the computation using hardware acceleration with GPUs. One specific line of research along these lines is finding privacy-preserving variants of GPUPtreeSHAP proposed in [MFH22].

## References

- [Act23] EU AI Act. European Union Artificial Intelligence Act. <https://artificialintelligenceact.eu/>, 2023.
- [aM23] Scale and Mamba. Scale and mamba. <https://github.com/KULeuven-COSIC/SCALE-MAMBA>, 2023.
- [BCD<sup>+</sup>23] M. Georgieva Belorgey, S. Carпов, K. Deforth, D. Jetchev, A. Sae-Tang, M. Vuille, N. Gama, J. Katz, I. Leontiadis, and M. Mohammadi. Manticore: A framework for efficient multiparty computation supporting real number and boolean arithmetic. *J. Cryptol.*, 36(3):31, 2023. doi:10.1007/s00145-023-09464-4.
- [BCG<sup>+</sup>18] Christina Boura, Ilaria Chillotti, Nicolas Gama, Dimitar Jetchev, Stanislav Peceny, and Alexander Petric. High-precision privacy-preserving real-valued function evaluation. In *Financial Cryptography and Data Security*, pages 183–202, Berlin, Heidelberg, 2018. Springer Berlin Heidelberg. doi:10.1007/978-3-662-58387-6\_10.
- [BIS<sup>+</sup>22] A. Bogdanova, A. Imakura, T. Sakurai, T. Fujii, T. Sakamoto, and H. Abe. Achieving transparency in distributed machine learning with explainable data collaboration. *CoRR*, abs/2212.03373, 2022. arXiv:2212.03373, doi:10.48550/arXiv.2212.03373.
- [BLW08] D. Bogdanov, S. Laur, and J. Willemsen. Sharemind: A framework for fast privacy-preserving computations. In *European Symposium on Research in Computer Security*, pages 192–206. Springer, 2008. doi:10.1007/978-3-540-88313-5\_13.
- [CDE<sup>+</sup>18] R. Cramer, I. Damgård, D. Escudero, P. Scholl, and C. Xing. SPDZ<sub>2k</sub>: Efficient mpc mod  $2^k$  for dishonest majority. In *Advances in Cryptology – CRYPTO 2018*, pages 769–798, 2018. doi:10.1007/978-3-319-96881-0\_26.

- [DDG<sup>+</sup>22] K. Deforth, M. Desgroseilliers, N. Gama, M. Georgieva, D. Jetchev, and M. Vuille. XORBoost: Tree boosting in the multiparty computation setting. *Proc. Priv. Enhancing Technol.*, 2022(4):66–85, 2022. doi:10.56553/popets-2022-0099.
- [DH12] W.J. Dally and R.C. Harting. *Digital Design: A Systems Approach*. Cambridge University Press, 2012. URL: <https://books.google.ch/books?id=WLoHOGOMhbIC>.
- [DPSZ12] I. Damgård, V. Pastro, N. Smart, and S. Zakarias. Multiparty computation from somewhat homomorphic encryption. In *Annual Cryptology Conference*, pages 643–662. Springer, 2012. doi:10.1007/978-3-642-32009-5\_38.
- [EDPS23] European Data Protection Supervisor. EDPS TechDispatch on Explainable AI. [www.edps.europa.eu/system/files/2023-11/23-11-16\\_techdispatch\\_xai\\_en.pdf](http://www.edps.europa.eu/system/files/2023-11/23-11-16_techdispatch_xai_en.pdf), 2023.
- [EGK<sup>+</sup>20] D. Escudero, S. Ghosh, M. Keller, R. Rachuri, and P. Scholl. Improved primitives for MPC over mixed arithmetic-binary circuits. In *40th Annual International Cryptology Conference, CRYPTO*, volume 12171 of *Lecture Notes in Computer Science*, pages 823–852, 2020. doi:10.1007/978-3-030-56880-1\_29.
- [JSTV23] D. Jetchev, A. Sae-Tang, and M. Vuille. Balancing privacy and explainable ai in semiconductor manufacturing. <https://inpher.io/blog/xai-semiconductor-manufacturing/>, 2023.
- [Kel20] M. Keller. MP-SPDZ: A versatile framework for multi-party computation. In *CCS '20: 2020 ACM SIGSAC Conference on Computer and Communications Security*, pages 1575–1590, 2020. doi:10.1145/3372297.3417872.
- [KOS16] M. Keller, E. Orsini, and P. Scholl. Mascot: faster malicious arithmetic secure computation with oblivious transfer. In *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security*, pages 830–842, 2016. doi:10.1145/2976749.2978357.
- [KPR18] M. Keller, V. Pastro, and D. Rotaru. Overdrive: Making SPDZ great again. In *EUROCRYPT 2018*, volume 10822 of *Lecture Notes in Computer Science*, pages 158–189, 2018. doi:10.1007/978-3-319-78372-7\_6.
- [LEC<sup>+</sup>20] S. Lundberg, G. Erion, H. Chen, A. DeGrave, J. Prutkin, B. Nair, R. Katz, J. Himmelfarb, N. Bansal, and S.-I. Lee. From local explanations to global understanding with explainable AI for trees. *Nat. Mach. Intell.*, 2(1):56–67, 2020. doi:10.1038/s42256-019-0138-9.
- [LEL18] S. Lundberg, G. Erion, and S. Lee. Consistent individualized feature attribution for tree ensembles. *CoRR*, abs/1802.03888, 2018.
- [LL17] S. Lundberg and S.-I. Lee. A unified approach to interpreting model predictions. In *Advances in Neural Information Processing Systems*, pages 4765–4774, 2017.
- [MFH22] Rory Mitchell, Eibe Frank, and Geoffrey Holmes. Gputreeshap: massively parallel exact calculation of shap scores for tree ensembles. *PeerJ Computer Science*, 8:e880, 04 2022. doi:10.7717/peerj-cs.880.
- [Mol22] C. Molnar. *Interpretable Machine Learning*. 2 edition, 2022. URL: <https://christophm.github.io/interpretable-ml-book>.

- [MR18] P. Mohassel and P. Rindal. *Aby<sup>3</sup>*: A mixed protocol framework for machine learning. In David Lie, Mohammad Mannan, Michael Backes, and XiaoFeng Wang, editors, *Proceedings of the 2018 ACM SIGSAC Conference on Computer and Communications Security, CCS 2018, Toronto, ON, Canada, October 15-19, 2018*, pages 35–52. ACM, 2018. doi:[10.1145/3243734.3243760](https://doi.org/10.1145/3243734.3243760).
- [MZ17] P. Mohassel and Y. Zhang. *SecureML*: A system for scalable privacy-preserving machine learning. In *2017 IEEE Symposium on Security and Privacy, SP 2017, San Jose, CA, USA, May 22-26, 2017*, pages 19–38. IEEE Computer Society, 2017. doi:[10.1109/SP.2017.12](https://doi.org/10.1109/SP.2017.12).
- [PSSY21] A. Patra, T. Schneider, A. Suresh, and H. Yalame. *Aby2. 0*: Improved mixed-protocol secure two-party computation. In *30th USENIX Security Symposium*, 2021.
- [RSG16] M. Túlio Ribeiro, S. Singh, and C. Guestrin. "why should I trust you?": Explaining the predictions of any classifier. In *Proceedings of the Demonstrations Session, NAACL HLT 2016, The 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, San Diego California, USA, June 12-17, 2016*, pages 97–101. The Association for Computational Linguistics, 2016. doi:[10.18653/v1/N16-3020](https://doi.org/10.18653/v1/N16-3020).
- [Sha53] L. Shapley. A value for n-person games. *Contributions to the Theory of Games*, 2(28):307–317, 1953. doi:[10.1515/9781400881970-018](https://doi.org/10.1515/9781400881970-018).
- [SK14] E. Strumbelj and I. Kononenko. Explaining prediction models and individual predictions with feature contributions. *Knowl. Inf. Syst.*, 41(3):647–665, 2014. doi:[10.1007/s10115-013-0679-x](https://doi.org/10.1007/s10115-013-0679-x).
- [TLL<sup>+</sup>22] Z. Tian, J. Liu, J. Li, X. Cao, R. Jia, and K. Ren. Private data valuation and fair payment in data marketplaces. *CoRR*, abs/2210.08723, 2022. arXiv:[2210.08723](https://arxiv.org/abs/2210.08723), doi:[10.48550/arXiv.2210.08723](https://doi.org/10.48550/arXiv.2210.08723).
- [Wan19] Guan Wang. Interpret federated learning with shapley values, 05 2019. doi:[10.48550/arXiv.1905.04519](https://doi.org/10.48550/arXiv.1905.04519).
- [WAYS22] Lauren Watson, Rayna Andreeva, Hao-Tsung Yang, and Rik Sarkar. Differentially private shapley values for data evaluation, 06 2022. doi:[10.48550/arXiv.2206.00511](https://doi.org/10.48550/arXiv.2206.00511).
- [WGC19] S. Wagh, D. Gupta, and N. Chandran. *SecureNN*: 3-party secure computation for neural network training. *Proceedings on Privacy Enhancing Technologies*, 2019(3):26–49, 2019. doi:[10.2478/popets-2019-0035](https://doi.org/10.2478/popets-2019-0035).
- [Yan22] Jilei Yang. Fast treeshap: Accelerating shap value computation for trees, 2022. arXiv:[2109.09847](https://arxiv.org/abs/2109.09847), doi:[10.48550/arXiv.2109.09847](https://doi.org/10.48550/arXiv.2109.09847).

## A Table of Notations

**Table 5:** Mathematical notation related to XorSHAP

Notation specific to SHAP and TreeSHAP:

$f_S(x)$	worth function for sample $x$ and subset $S$ (§2.1)
$\phi_i(x)$	Shapley value for sample $x$ and feature $i$ (§2.1)
$\mathcal{P}_\ell$	path leading to leaf $\ell$ (§2.2)
$\mathcal{J}_\ell$	multiset of the feature indices for split nodes in $\mathcal{P}_\ell$ (§2.2)
$\mathcal{F}_\ell$	distinct feature indices for split nodes in $\mathcal{P}_\ell$ (size at most $D$ )
$c_n$	cover of node $n$ in a tree (§2.2)
$W_{\ell,S}(x)$	proportion of paths that flow down to leaf node $\ell$ (§2.2)

Notation specific to the XorSHAP data-independent algorithm:

$P_\ell$	$D \times M$ binary matrix of feature selector vectors for the path $\mathcal{P}_\ell$ (§3.1)
$F_\ell$	characteristic vector of $\mathcal{F}_\ell$ (of size $M$ )
$\sigma_{F_\ell}$	oblivious sorting permutation for $F_\ell$ (§3.2)
$\mathcal{R}_\ell$	superset of $\mathcal{F}_\ell$ of the last $D$ indices after applying $\sigma_{F_\ell}$ (§3.2)
$R_\ell$	characteristic vector of $\mathcal{R}_\ell$ (of size $D$ )
$\mathbf{b}_j$	feature selector vector for feature $j$
$\mathbf{o}_\ell$	binary vector of size $D$ (§3.1)
$\mathbf{z}_\ell$	column vector of size $D$ with cover ratios (§3.1)
$\mathbf{t}_\ell$	vector of size $M$ used to compute Shapley values (§3.2)
$C_S$	Shapley weight (§3.2.4)

Other notation used in the document:

$b ? x : y$	Ternary conditional operator (ternary if), obviously computed as $b ? x : y = y + b \cdot (x - y)$
$n$	a node of a binary decision tree
$n_{\text{left}}$	the left child of a split node $n$
$n_{\text{right}}$	the right child of a split node $n$
$\ell$	a leaf node of a binary decision tree
$w_\ell$	the weight of a leaf node of a binary decision tree