

# Privacy-Preserving Contrastive Explanations with Local Foil Trees

Thijs Veugen<sup>1,2</sup>[0000-0002-9898-4698], Bart Kamphorst<sup>1</sup>[0000-0002-9490-5841],  
and Michiel Marcus<sup>1</sup>[0000-0003-0936-2289]

<sup>1</sup> TNO, The Hague, The Netherlands

`thijs.veugen,bart.kamphorst,michiel.marcus@tno.nl`

<http://www.tno.nl>

<sup>2</sup> CWI, Amsterdam, The Netherlands

<http://www.cwi.nl>

**Abstract.** We present the first algorithm that combines privacy-preserving technologies and state-of-the-art explainable AI to enable privacy-friendly explanations of black-box AI models. We provide a secure algorithm for contrastive explanations of black-box machine learning models that securely trains and uses local foil trees. Our work shows that the quality of these explanations can be upheld whilst ensuring the privacy of both the training data, and the model itself.

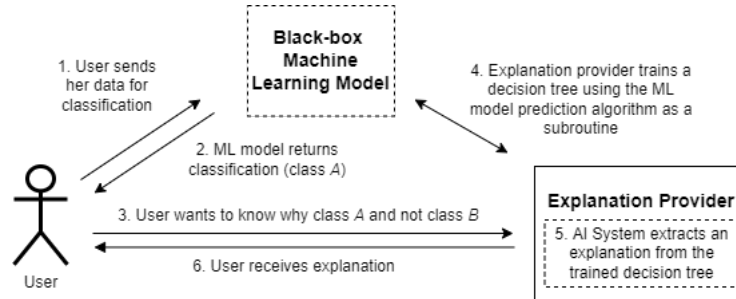
**Keywords:** Explainable AI · secure multi-party computation · decision tree · foil tree

## 1 Introduction

The field of explainable AI focuses on improving the interpretability of machine learning model behaviour. In recent years, exciting developments took place in this area, such as the emergence of the LIME [15] and SHAP [13] algorithms, which have become popular. These algorithms take a data point and its classification according to a trained machine learning model, and provide an explanation for the classification by analyzing the importance of each feature for that specific classification. This is interesting for a researcher, but a layman using the AI system is unlikely to understand the reasoning of the machine learning model.

Instead, Van der Waa et al. [20] created an algorithm called *local foil trees* that explains why someone was classified as class  $A$  instead of another class  $B$ , by providing a set of decisions rules that need to apply for that point to be classified as class  $B$ . This provides an increased understanding of the AI system [19], which can for instance be used to infer what can be done to change the classification. This is particularly relevant to decision support systems, where the AI system should provide advice to the user. An example could be that the AI system advises a user to have lower blood pressure and higher body weight, in order to go from high risk of a certain illness to a lower risk.

Our work focuses on creating a secure algorithm that provides the same functionality as the local foil tree algorithm in a setting where the black-box



**Fig. 1.** Overview of steps and interactions in the local foil tree algorithm.

machine learning model needs to remain secret to protect the confidentiality of the machine learning model and the training data. Before we explain why this assumption is realistic, we provide a rough overview of the algorithm and interactions in the local foil tree algorithm.

As shown in Figure 1, the user first submits her data to the machine learning model to retrieve a classification. The user then wants to know why she was classified as class *A* and not as class *B*. To create an explanation for this, the explanation provider trains a decision tree and uses the machine learning model as a black-box subroutine within that process. This decision tree is then used to generate an explanation.

In practice, we often see that it can be very valuable to train machine learning models on personal data, for example in the medical domain to prevent diseases [21], or to detect possible money laundering [16]. Due to the sensitive nature of personal data, however, it is challenging for organisations to share and combine data. Legal frameworks like the General Data Protection Regulation<sup>3</sup> (GDPR) and the Health Insurance Portability and Accountability Act<sup>4</sup> (HIPAA) further restrict the usage and exchange of personal data.

In order to avoid violating privacy when we want to use personal data as training data for a machine learning algorithm, it is possible to apply cryptographic techniques to securely train the machine learning model, which results in a hidden model [21,12,22]. This ensures that the privacy of the personal data is preserved while it is used to train the model. In order to enable explainable AI with the hidden model, we could simply reveal the model and apply, e.g., the original local foil tree algorithm. However, there are various reasons why it could be undesirable to reveal the trained model. Firstly, if one or more organisations involved have a commercial interest in the machine learning model, the model could be used in ways that were not originally intended. Keeping the model secret then ensures control of model usage. Secondly, sensitive data is

<sup>3</sup> <https://gdpr-info.eu>

<sup>4</sup> <https://www.govinfo.gov/content/pkg/PLAW-104publ191/pdf/PLAW-104publ191.pdf>

used to train the machine learning model and recent research has shown that it is feasible to reconstruct training data from a trained model [9,23,25]. The whole reason to securely train the model is to avoid leaking sensitive data, but if the machine learning model is known, it is still possible that sensitive data is leaked when such reconstruction attacks are used. In these cases, we should therefore assume that the model stays hidden to protect the confidentiality of the machine learning model and the training data.

This poses a new challenge for black-box explainable AI. In step 2 of Figure 1 the classification  $A$  can be revealed to the user without problems, but it is unclear how steps 4 and 5 from Figure 1 would work when the model is hidden. There is a variety of cryptographic techniques that can be used to securely train models. When multiple organisations are involved, common techniques are secret sharing [5] and homomorphic encryption [14]. In this work, we address the aforementioned challenge and provide an algorithm that can produce contrastive explanations when the model is either secret shared, or homomorphically encrypted. Practically, this means that the explanation provider, as shown in Figure 1, does not have the model locally, but that it is owned by a different party or even co-owned by multiple parties. The arrows in the figure then imply that communication needs to happen with the parties that (jointly) own the model.

An additional challenge comes from the fact that explainable AI works best when rule-based explanations, as provided through the local foil tree algorithm, are accompanied by an example-based explanation, such as a data point that is similar to the user, but is classified as class  $B$  instead of  $A$  [19]. The use of a data point (having class  $B$ ) from the sensitive training data would violate privacy in the worst way possible. As we will discuss in section 3, we address this challenge using synthetic data.

In summary, we present a privacy-preserving solution to explain an AI system, consisting of:

- A cryptographic protocol to securely train a binary decision tree when the target variable is hidden;
- An algorithm to securely generate synthetic data based on numeric sensitive data;
- A cryptographic protocol to extract a rule-based explanation from a hidden foil tree, and construct an example data point for it.

The target audience for this work is twofold. On the one hand, our work is relevant for data scientists who want to provide explainable, data-driven insights using sensitive (decentralized) data. It gives access to new sources of data without violating privacy when explainability is essential. On the other hand, our work provides a new tool for cryptographers to improve the interpretability of securely trained machine learning models which have applications in the medical and financial domain.

In the remainder of this introduction, we discuss related work and briefly introduce secure multi-party computation. In the sections following after, we

explain the local foil tree algorithm [20] and present a secure solution. Thereafter, we discuss the complexity of the proposed solution and share experimental results. Finally, we provide closing remarks in the conclusion.

### 1.1 Related work

Our solution is based on the local foil tree algorithm by Van der Waa et al. [20], for which we design a privacy-preserving solution based on MPC. There is related work in the area of securely training decision trees, but these results are never applied to challenges in explainable AI. As we will elaborate on further in section 3, we have a special setting where the feature values of the synthetic data to train the decision tree are not encrypted, but the classifications of these data points *are* encrypted. As far as we know, no training algorithm for such a setting has been proposed yet.

We mention the work of de Hoogh et al. [7], who present a secure variant of the well-known ID3 algorithm (with discrete variables). Their training data points remain hidden, whereas in our case, that is not necessary. Furthermore, as the number of children of an ID3 decision node reflects the number of categories of the chosen feature, the tree decision is not completely hidden. They implement their solution using Shamir sharing with VIFF, which is a predecessor of the MPyC [17] framework that we use.

A more recent paper on secure decision trees is by Abspoel et al. [1], who implement C4.5 and CART in the MP-SPDZ framework. We also consider CART since this leads to a binary tree, which does not reveal information on (the number of categories of) the feature chosen in a decision node. Abspoel et al. use both discrete and continuous variables, similar to our setting. However, since Abspoel et al. work with encrypted feature values, they need a lot of secure comparisons to determine the splitting thresholds.

In a similar approach, Adams et al. [2] scale the continuous features to a small domain to avoid the costly secure comparisons, at the expense of a potential drop in accuracy.

Only one article was found on privacy-preserving explainable AI. The work of [11] presents a new class of machine learning models that are interpretable and privacy-friendly with respect to the training data. Our work does not introduce new models, but provides an algorithm to improve the interpretability of existing complex models that have been securely trained on sensitive data.

### 1.2 Secure multi-party computation

We use secure multi-party computation (MPC) to protect secret data, such as the ML classification model, and its training data. MPC is a cryptographic tool to extract information from the joint data of multiple parties, without needing to share their private data with other parties. Introduced by Yao in 1982 [24], the field has developed fast, and various platforms are available now for arbitrary secure computations on secret data, such as addition, subtraction, multiplication and comparison. We use the MPyC platform [17] that uses Shamir secret sharing

$A$	Fact (class); classification of the user as indicated by the black-box.
$B$	Foil (class); target class for contrastive explanation to the user.
$\mathcal{B}$	Decision tree or, equivalently, foil tree.
$G_s$	Gini index for split $s \in \{1, \dots, \varsigma\}$ .
$\tilde{G}_s = N_s/D_s$	Adjusted Gini index for split $s \in \{1, \dots, \varsigma\}$ .
$k_A, k_B$	Index of classes $A$ and $B$ , respectively.
$K$	Number of classes.
$n$	Number of available synthetic data points in a particular node.
$N$	Number of synthetic data points $ \mathcal{X} $ .
$P$	Number of features per data point.
$\varsigma$	Number of splits $ \mathcal{S} $ .
$S_s = (p_s, t_s)$	Feature index $p_s \in \{1, \dots, m\}$ and threshold $t_s$ of split $S_s$ , $1 \leq s \leq \varsigma$ .
$\mathbf{x}_i, \mathbf{x}_U$	Vector $(x_{i,1}, \dots, x_{i,P})$ of feature values of synthetic data point $i$ . With subscript $U$ , it refers to the data point of the user.
$\mathcal{X}$	Set of all synthetic data points $\mathbf{x}_i$ , $i = 1, \dots, N$ .
$\mathbf{y}_i$	Indicator vector $(y_{i,1}, \dots, y_{i,K})$ of the class of data point $i$ as indicated by the black-box.
$\xi_i$	Bit that indicates whether data point $i$ is available (1) or unavailable (0) in the current node.

**Table 1.** Notation as used throughout the document. Some symbols are seen in the context of a certain point (node) within the decision tree, in which case they can be sub- or superscripted with  $l$  or  $r$  to denote the same variable in the left or right child node that originates from the current node.

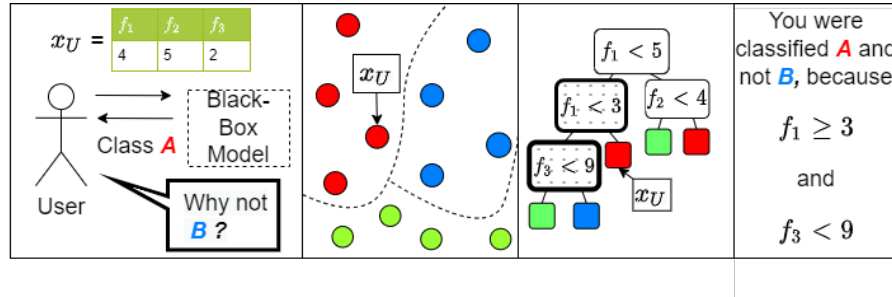
in the semi-honest model, where all parties are curious, but are assumed to follow the rules of the protocol.

Like many MPC platforms, MPyC follows the share-compute-reveal paradigm. Each party first uploads its inputs, by generating non-revealing shares for the other parties. When the inputs have been uploaded as secrets, the parties can then perform joint computations without learning the inputs. Finally, the eventually computed output is revealed to the entitled parties.

### 1.3 Notation

Due to the inherent complexity of both explainable AI and cryptographic protocols, we require many symbols in our presentation. These symbols are all introduced in the body of this paper; however, for the reader’s convenience we also summarize the most important symbols in Table 1.

Sets are displayed in curly font, e.g.  $\mathcal{X}$ , and vectors in bold font, e.g.  $\mathbf{x}_U$ . The vector  $\mathbf{e}_j$  represents the  $j$ -th elementary vector of appropriate, context-dependent length. The notation  $(x \geq y)$  is used to denote the Boolean result of the comparison  $x \geq y$ . Any symbol between square brackets  $[\cdot]$  represents a secret-shared version of that symbol. Finally, a reference to line  $y$  of Protocol  $x$  is formulated as line  $x.y$ .



**Fig. 2.** A visualisation of the different steps in the local foil tree algorithm to explain why data point  $x_U$  was classified as (fact) class  $A$  and not as (foil) class  $B$ . The different images depict classification retrieval and foil class selection, data preparation, decision tree training and determining relevant nodes, and explanation extraction.

---

### Protocol 1 Foil-tree based explanation

---

**Input:** Data point  $x_U$  that is classified as class  $A$ ; foil class  $B$

**Output:** Explanation why  $x_U$  was not classified as the foil class

- 1: Obtain a classification for the user ▷ cf. Section 3.1
  - 2: Prepare the synthetic data points for the foil tree ▷ 3.2
  - 3: Classify all synthetic data points through the black-box ▷ 3.3
  - 4: Train a decision tree ▷ 3.4
  - 5: Locate fact leaf (leaf node of  $x_U$ ) ▷ 3.5
  - 6: Determine the foil leaf (leaf node of class  $B$  closest to fact leaf) ▷ 3.6
  - 7: Determine the decision node at which the root-leaf paths of the fact and foil leaf split ▷ 3.7
  - 8: Construct the explanation (and provide example data point). ▷ 3.7
- 

## 2 Explainable AI with Local Foil Trees

In this section we present the local foil tree method of Van der Waa et al. [20] and discuss the challenges that arise, when the black-box classifier does not yield access to its training data, and provides classifications in secret-shared (or encrypted) form to the explanation provider.

We assume that we have black-box access to a classification model. If a user-supplied data point  $x_U$  is classified as some class  $A$ , our goal is to construct an explanation why  $x_U$  was not classified as another class  $B$ . The explanation will contain decision rules of the form that a certain feature of  $x_U$  is less (or greater) than a certain threshold value. An overview of the different steps is illustrated in Figure 2 and formalized in Protocol 1. Note that we deviate from Van der Waa et al. by providing an example data point in the final step. In each step of the protocol, we also refer to the section of our work where we present secure protocols for that step.

Both limitations that the model owner introduced help to better preserve the privacy of the training data, and to secure the model itself, but also hinder us

in generating an explanation. In particular, if we cannot have access to training data, we need to generate synthetic data that can be used to train a decision tree. Training a decision tree in itself is complicated by the fact that the classifications are hidden, which most notably implies that during the recursive procedure we need to securely keep track of the synthetic data samples that end up in each branch of the tree.

Information can also be revealed through the structure of the decision tree; in particular, it may disclose the splitting feature. For example, if a certain categorical feature can assume six values and a decision node splits into six new nodes, it is likely that this node represents that feature. For this reason we do not use the commonly-used ID3 or its successor C4.5 for training the decision tree. We instead generate a binary decision tree with the CART (Classification and Regression Trees) algorithm [3]. The CART algorithm greedily picks the best decision rule for a node. In case of classification trees, this materializes as the rule with the lowest Gini index. The Gini index measures the impurity, i.e. the probability of incorrectly classifying an observation, so the lowest Gini index results in the best class purity.

The result of the training procedure is a decision tree whose decision rules and leaf classification are secret-shared. As a consequence we need a secure protocol for determining the position of a foil data point, and all nodes that are relevant for the explanation. With help of the model owner(s), among all secret values in the process, only these nodes, and the user classification are revealed.

Compared to secure protocols for training decision trees on hidden data points with hidden classification, the fact that we use synthetic data also has some benefits. First, since the (synthetic) data points are known, we can still access their features directly, improving the efficiency of the protocol. Second, since we already trained the decision tree on synthetic data, we can also supplement our explanation with a synthetic data point, and thereby increase user acceptance [19].

### 3 Secure solution

In this section we describe the secure version of the local foil tree algorithm, which reveals negligible information about the sensitive training data and black-box model. In the rest of this work, we will refer to *training data* when we talk about the data used to train the black-box machine learning model and to *synthetic data* when we refer to the synthetically generated data that we use to train the foil tree.

The secure protocol generates  $N$  synthetic data points  $\mathbf{x}_i$ ,  $i = 1, \dots, N$ , with  $P$  features that each can be categorical, or continuous. To increase the efficiency of the secure solution, we make use of one-hot or integer encoding to represent categorical values. We assume that the class  $k \in \{1, \dots, K\}$  of data point  $\mathbf{x}_i$  is represented by a secret binary indicator vector  $[\mathbf{y}_i] = ([y_{i,1}], \dots, [y_{i,K}])$ , such that  $y_{i,k} = 1$ , if data point  $\mathbf{x}_i$  is classified as class  $k$  by the black-box, and  $y_{i,k} = 0$ , otherwise.

During the decision training, we maintain an indicator vector  $\xi$  of length  $N$ , such that  $\xi_i = 1$ , if and only if, the  $i$ -th synthetic data point is still present in this branch.

### 3.1 Classify user data

We assume that the user is allowed to learn the black-box classification of her own data point  $\mathbf{x}_U$ , so this step is trivial. Without loss of generality, we assume that the user received classification  $A$ .

### 3.2 Generating synthetic data

Van der Waa et al. [20] mention that synthetic data could be used to train the local foil trees, and suggest using normal distributions. In this section, we apply that suggestion and provide a concrete algorithm for generating a local data set around the data point for which an explanation is being generated.

We first take a step back and list what requirements the synthetic data should adhere to:

1. The synthetic data should reveal negligible information about the features of the training data.
2. The synthetic data should be local, in the sense that all data points are close to  $\mathbf{x}_U$ , the data point to be explained.
3. The synthetic data should be realistic, such that they can be used in an explanation and still make sense in that context.

State-of-the-art synthetic data generation algorithms, such as GAN [10] and SMOTE [4] can generate very realistic data, but they need more than one data point to work, so we cannot apply them to the single data point to be explained. One could devise a secure algorithm for GAN or SMOTE and securely apply it to the sensitive data, but this would affect the efficiency of our solution. In this article, we pursue the simpler approach that was suggested by Van der Waa et al.

Ideally, one would securely calculate some statistics of the sensitive training data for the black-box model and reveal these statistics. Based on these statistics, one could generate a synthetic data set by sampling from an appropriate distribution. Our implementation securely computes the mean and variance of every feature in the training data and samples synthetic data points from a truncated normal distribution. The reason for truncating is two-fold: first, it allows us to sample close to the user’s data point  $\mathbf{x}_U$ , and second, features may not assume values on the entire real line. Using a truncated normal distribution allows us to generate slightly more realistic data that is similar to  $\mathbf{x}_U$ . The details are presented in Protocol 2.

To generate more realistic data, one could also incorporate correlation between features, or go even further and sample from distributions that better represent the training data than a normal distribution.



---

**Protocol 2** Synthetic data generation.
 

---

**Input:** Encrypted black-box training set  $[\tilde{\mathcal{X}}]$ , integer  $N$ , target data point  $\mathbf{x}_U$   
**Output:** Synthetic data set  $\mathcal{X}$  with cardinality  $N$

- 1: **for**  $p = 1, \dots, P$  **do** ▷ Compute mean and variance of feature  $p$
- 2:      $[\mu_p] \leftarrow 1/|\tilde{\mathcal{X}}| \sum_{\tilde{\mathbf{x}}_U \in \tilde{\mathcal{X}}} \tilde{x}_p$ ;      $[\sigma_p^2] \leftarrow 1/|\tilde{\mathcal{X}}| \sum_{\tilde{\mathbf{x}}_U \in \tilde{\mathcal{X}}} (\tilde{x}_p - [\mu_p])^2$
- 3: **end for**
- 4: Reveal  $\boldsymbol{\mu}$  and  $\boldsymbol{\sigma}^2$
- 5:  $\mathcal{X} \leftarrow \emptyset$
- 6: **for**  $i = 1, \dots, N$  **do**
- 7:     **for**  $p = 1, \dots, P$  **do**
- 8:         **repeat** Draw  $x_{i,p}$  from  $\mathcal{N}(\mu_p, \sigma_p^2)$
- 9:         **until**  $x_{i,p} \in [x_p - 3\sigma_p, x_p + 3\sigma_p]$
- 10:     **end for**
- 11:      $\mathcal{X} \leftarrow \mathcal{X} \cup \mathbf{x}_i$
- 12: **end for**
- 13: Return  $\mathcal{X}$

---

In our experiments, we noticed that an interval of  $[x_p - 3\sigma_p, x_p + 3\sigma_p]$  generally yielded a synthetic data set that was still close to  $\mathbf{x}_U$ , but also provided a variety of classifications for the data points. A smaller interval (for example of size  $2\sigma_p$ ) often resulted in a data set for which the distribution of classifications was quite unbalanced. The foil class might then not be present in the foil tree, breaking the algorithm. Larger intervals would result in data points that are not local anymore, and therefore yield a less accurate decision tree.

### 3.3 Classify synthetic data

All synthetic training data point  $\mathbf{x}_i$  can now be classified securely by the model owner(s). This results in secret-shared classification vectors  $[\mathbf{y}_i]$ . The secure computation depends on the model, and is beyond our scope.

### 3.4 Training a decision tree

In this section, we explain the secure CART algorithm that we use to train a secure decision tree, which is formalized in Protocol 3. The inputs to this algorithm are:

1.  $\mathcal{X}$ : a set of synthetic data points.
2.  $\mathcal{S}$ : a set of splits to use in the algorithm. Each split  $S_s \in \mathcal{S}, s = 1, \dots, \zeta$  is characterized by a pair  $(p_s, t_s)$  that indicates that the feature with index  $p_s$  is at least  $t_s$ .
3.  $\tau$ : the convergence fraction used in the stopping criterion.
4.  $[\boldsymbol{\xi}]$ : a secret binary availability vector of size  $N$ . Here  $\xi_i$  equals 1, if the  $i$ -th synthetic data point is available, and 0, otherwise.

We start with an empty tree and all training data points are marked as available. First, the stopping criterion uses the number of elements of the most common class (line 3.8), and the total number of elements in the availability vector (line 3.7). The stopping criterion from line 3.10 is securely computed by  $1 - (1 - [(n \leq \tau \cdot N)] \cdot (1 - [(n = n_{k^*}]])$ , and consequently revealed.

If the stopping condition is met, i.e., equal to one, a leaf node with the secret-shared indicator vector of the most common class is generated. In order to facilitate the efficient extraction of a foil data point as mentioned at the start of section 3, we also store the availability vector  $\xi$  in this leaf node. How this indicator vector is used to securely generate a foil data point is discussed in section 3.8.

If the stopping criterion is not met, a decision node is created by computing the best split (lines 3.13–19) using the adjusted Gini indices of each split in  $\mathcal{S}$ . We elaborate on computing the adjusted Gini index (lines 3.13–15) later on in this section.

After determining the optimal split, an availability vector is constructed for each child based on this split in lines 3.21–22. For the left child, this is done using the availability vector  $[\xi]$ , the indicator vector indicating the feature  $p_{s^*}$  of the best split  $e_{p_{s^*}}$  and the threshold  $t_{s^*}$  of the best split as explained in protocol 4. The resulting availability vector  $[\xi^l]$  has a [1] in index  $i$ , if  $x_{i,p_{s^*}} \leq t_{s^*}$ . The entry-wise difference with  $[\xi]$  then gives the availability vector for the right child. The CART algorithm is then called recursively with the new availability vectors to generate the children of the decision node.

In protocol 3, we use two yet unexplained subroutines, namely **max** and **find**. The **max** subroutine securely computes the maximum value in a list using secure comparisons. Thereafter, the **find** subroutine finds the location of the maximum computed by **max** in the list that was input to **max**, which is returned as a secret-shared indicator vector indicating this location. The functions **max** and **find** are already implemented in MPyC. However, since we always use the two in junction, we implemented a slight variation that is presented in Appendix A.

### Compute the Gini index for each possible split

We aim to find the split  $S^* := S_{s^*} = (p_{s^*}, t_{s^*})$  with highest class purity, which is equivalent to the lowest Gini index  $G_s$ . As such, we first need to compute the Gini index for all splits. The Gini index of a split is the weighted sum of the Gini value of the two sets that are induced by the split,

$$G_s = g_s^l \cdot \frac{n^{s,l}}{n} + g_s^r \cdot \frac{n^{s,r}}{n}. \quad (1)$$

Here,  $n$  is again the number of available data points in the current node,  $n^{s,l}$  is the number of available data points in the left set that is induced by split  $S_s$ , and  $g_s^l$  is the Gini value of the left set that is induced by split  $S_s$ ,

$$g_s^l := 1 - \sum_{k=1}^K \left( \frac{n_k^{s,l}}{n^{s,l}} \right)^2, \quad (2)$$

---

**Protocol 3 cart**

Secure CART training of a binary decision tree.

---

**Input:** Training set  $\mathcal{X}$ , split set  $\mathcal{S}$ , convergence parameter  $\tau \in [0, 1]$ , secret-shared binary availability vector  $[\xi]$   
**Output:** Decision tree  $\mathcal{B}$

- 1:  $\mathcal{B} \leftarrow \emptyset$
- 2:  $N \leftarrow |\mathcal{X}|$
- 3: **while**  $\mathcal{B}$  is not fully constructed **do**
- 4:     **for**  $k = 1, \dots, K$  **do**
- 5:          $[n_k] \leftarrow \sum_{i=1}^N [y_{i,k}] \cdot [\xi_i]$  ▷ nr available data points per class
- 6:     **end for**
- 7:      $[n] \leftarrow \sum_{k=1}^K [n_k]$  ▷ nr available data points
- 8:      $[n_{k^*}] \leftarrow \max([n_1], \dots, [n_K])$
- 9:      $[e_{k^*}] \leftarrow \mathbf{find}([n_{k^*}], ([n_1], \dots, [n_K]))$  ▷ indicates most common class
- 10:     **if**  $[(n \leq \tau \cdot N)]$  or  $[(n = n_{k^*})]$  **then** ▷ branch fully constructed
- 11:         Extend  $\mathcal{B}$  with leaf node with class indicator  $[e_{k^*}]$
- 12:     **else** ▷ branch splits
- 13:         **for**  $s = 1, \dots, \zeta$  **do**
- 14:              $[G_s] \leftarrow \mathbf{adjusted\_gini}(S_s)$
- 15:         **end for**
- 16:          $[G_{s^*}] \leftarrow \max([G])$
- 17:          $[e_{k^*}] \leftarrow \mathbf{find}([G_{s^*}], [G])$  ▷ indicates best split
- 18:          $[p_{s^*}] \leftarrow \sum_{s=1}^{\zeta} [e_{s^*,s}] \cdot p_s$  ▷ feature of optimal split
- 19:          $[t_{s^*}] \leftarrow \sum_{s=1}^{\zeta} [e_{s^*,s}] \cdot t_s$  ▷ threshold of optimal split
- 20:          $b \leftarrow$  decision node that corresponds with split  $([p_{s^*}], [t_{s^*}])$
- 21:          $[\xi^l] \leftarrow \mathbf{left\_child\_availability}(\mathcal{X}, [\mathbf{x}i], [p^*], [t_{s^*}])$
- 22:          $[\xi^r] \leftarrow [\xi] - [\xi^l]$
- 23:         Extend  $b$  to the left with result of  $\mathbf{cart}(\mathcal{X}, \mathcal{S}, \tau, [\xi^l])$
- 24:         Extend  $b$  to the right with the result of  $\mathbf{cart}(\mathcal{X}, \mathcal{S}, \tau, [\xi^r])$
- 25:         Extend  $\mathcal{B}$  with  $b$
- 26:     **end if**
- 27: **end while**
- 28: Return  $\mathcal{B}$

---

where  $n_k^l$  denote the number of available data points in the left node with class  $k$ . The symbols  $n^{s,l}$ ,  $n_k^{s,l}$  and  $g_s^r$  are defined analogously for the right set. For notation convenience, justified as the upcoming derivations concern a fixed index  $s$ , we drop the superscripts  $s$  from the symbol  $n$ .

We now derive a more convenient expression for the Gini index. Substituting expression (2) into (1) and rewriting yields

$$G_s = \frac{n^l + n^r}{n} - \frac{n^r \sum_{k=1}^K (n_k^l)^2 + n^l \sum_{k=1}^K (n_k^r)^2}{n \cdot n^l \cdot n^r}. \quad (3)$$

---

**Protocol 4** left\_child\_availability

---

Indicate the data points that flow into the left child.

**Input:** Synthetic data set  $\mathcal{X}$ , availability vector  $[\xi]$  for the current node, feature indicator vector  $[e_{p_s}]$ , threshold  $[t_s]$ **Output:** Availability vector  $[\xi^l]$  for the left child

- 1: **for**  $i=1, \dots, N$  **do**
  - 2:    $[x_{i,p_s}] \leftarrow \sum_{p=1}^P [e_{p_s,p}] \cdot x_{i,p}$
  - 3:    $[\delta_i] \leftarrow [(x_{i,p_s} \leq t_s)]$
  - 4:    $[\xi_i^l] \leftarrow [\xi_i] \cdot [\delta_i]$
  - 5: **end for**
  - 6: Return  $[\xi^l]$
- 

---

**Protocol 5** adjusted\_gini

---

Compute the adjusted Gini index of a split.

**Input:** Synthetic data set  $\mathcal{X}$ , vector of available transactions  $\xi$ , split  $(p_s, t_s) = S_s \in \mathcal{S}$ **Output:** Encrypted numerator and denominator of adjusted Gini index  $[\tilde{G}_s] = [N_s]/[D_s]$ 

- 1: **for**  $i=1, \dots, N$  **do**
  - 2:    $\delta_i \leftarrow (x_{i,p_s} \leq t_s)$   $\triangleright$  1 if data point meets split criterion, else 0
  - 3: **end for**
  - 4:  $[n] \leftarrow \sum_{i=1}^N [\xi_i]$ ,  $[n^l] \leftarrow \sum_{i=1}^N \delta_i \cdot [\xi_i]$ ,  $[n^r] \leftarrow [n] - [n^l]$
  - 5:  $[n_k] \leftarrow \sum_{i=1}^N [y_{i,k}] \cdot [\xi_i]$ ,  $[n_k^l] \leftarrow \sum_{i=1}^N \delta_i \cdot [y_{i,k}] \cdot [\xi_i]$ ,  $[n_k^r] \leftarrow [n_k] - [n_k^l]$
  - 6: Return  $[N_s] \leftarrow [n^r] \sum_{k=1}^K ([n_k^l])^2 + [n^l] \sum_{k=1}^K ([n_k^r])^2$  and  $[D_s] \leftarrow [n^l] \cdot [n^r]$
- 

Now, since  $n = n^l + n^r$  is independent of the split, minimizing the Gini index over all possible splits is equivalent to *maximising* the *adjusted Gini index*  $\tilde{G}_s$ ,

$$\tilde{G}_s = \frac{n^r \sum_{k=1}^K (n_k^l)^2 + n^l \sum_{k=1}^K (n_k^r)^2}{n^l \cdot n^r} =: \frac{N_s}{D_s}. \quad (4)$$

We represent  $\tilde{G}_s$  as a rational number to avoid an expensive secure (integer) division. Both the numerator  $N_s$  and denominator  $D_s$  are non-zero, if the split  $S_s$  separates the available data points, e.g., the split induces at least one available data point in each set. Otherwise, either  $n^l = 0$  or  $n^r = 0$ , and  $N_s = D_s = 0$ , such that  $\tilde{G}_s$  is not properly defined. In line 3.16 one could naively let  $\max$  evaluate ( $\tilde{G}_1 < \tilde{G}_2$ ) by computing ( $N_1 D_2 < N_2 D_1$ ). However, this may yield an undesired outcome if one of the denominators equals zero. Appendix B presents two possible modifications that handle this situation.

Protocol 5 shows how the adjusted Gini index can be computed securely. Observe that  $[n]$  and  $[n_k]$  were already computed for the CART stopping criterion, so they come for free. The computation  $[n_k^l]$  can be implemented efficiently as a secure inner product. The computations of  $n$ ,  $n^l$ ,  $n^r$ , and  $n_k^r$  don't require any additional communication. Because the total number of possible splits  $N \cdot P$  is much larger than the number  $N$  of data points, it makes sense to precompute

$[y_{i,k}] \cdot [\xi_i]$  for each  $i$  and  $k$ , such that the computation of  $n_k^l$  for each split requires no additional communication.

### Convergence

In theory, it is possible that, at some point during training, the CART algorithm has not met the stopping criterion yet, and has no splits available that actually separate the set of available data points. In this case the algorithm keeps adding useless decision nodes and does not make any actual progress. To prevent ending up in this situation, we can detect it by revealing ( $D_{s^*} = 0$ ), and take appropriate action. Also, a maximum number of nodes, or a maximum depth, can be set.

### 3.5 Locate the fact leaf

Once the decision tree has been constructed, we need to find the leaf that contains the fact  $\mathbf{x}_U$ . As the fact leaf will be revealed, the path from the root to the fact leaf will be revealed as well. Therefore, we can traverse the decision tree from the root downwards and reveal each node decision. The decision for data point  $\mathbf{x}_U$  at a given node can be computed similarly to Protocol 4. First, the feature value that is relevant for the current decision node is determined through  $[x_{U,p_{s^*}}] = \sum_{i=1}^P [e_{p_{s^*},i}] \cdot [x_{i,p}]$ . Second, the secure comparison  $[(x_{U,p_{s^*}} \leq t_{s^*})]$  is performed and revealed. The result directly indicates the next decision node that needs to be evaluated. This process is repeated until a leaf is encountered: the fact leaf.

### 3.6 Locate the foil leaf

Since we know the fact leaf and the structure of the decision tree, we can create an ordered list of all tree leaves, starting with the closest leaf and ending with the farthestmost leaf. We can traverse this list and find the first leaf that is classified as class  $B$ , without revealing the classes, but only whether they equal  $B$  or not, i.e. by revealing the Boolean  $[(e_{k^*,k_B} = 1)]$  for every leaf. This does not require any extra computations, as these vectors have already been computed and stored during the training algorithm. We use the number of steps between nodes within the decision tree as our distance metric, but as Van der Waa et al. [20] note, there are more advanced options.

### 3.7 Construct the explanation

Once the fact leaf and the foil leaf have been determined, the lowest common node can be found without any secure computations, since the structure of the decision tree is known. We traverse the decision tree from this lowest common node to the foil leaf and reveal the feature and threshold for each of the nodes on that path (the nodes with a thick border and dotted background in Figure 2). For each rule, we determine whether it applies to  $\mathbf{x}_U$ . For instance, if a rule says

**Protocol 6** retrieve\_foil

Retrieve foil data point

**Input:** Availability vector  $[\xi]$  of the foil leaf, class index  $k_B$ **Output:** Foil data point  $\mathbf{s}$ 

- 
- 1:  $[\varepsilon] \leftarrow [0]$   $\triangleright$  flips to [1] when a foil data point is found
  - 2: **for**  $i = 1, \dots, n$  **do**
  - 3:      $[\delta_i] \leftarrow (1 - [\varepsilon]) \cdot [\xi_i] \cdot [y_{i,k_B}]$
  - 4:      $[\varepsilon] \leftarrow [\varepsilon] + [\delta_i]$
  - 5: **end for**
  - 6: **for**  $p = 1, \dots, P$  **do**
  - 7:      $[s_p] \leftarrow \sum_{i=1}^N [\delta_i] \cdot [x_{i,p}]$
  - 8: **end for**
  - 9: Reveal  $\mathbf{s}$  to the user
- 

that  $x_{U,i} \geq 3$  and indeed  $\mathbf{x}_U$  satisfies this rule, then it is not relevant for the explanation.

After this filter is applied, we combine the remaining rules where applicable. For example, if one rule requires  $x_{U,i} \geq 3$  and another rule requires  $x_{U,i} \geq 4$ , we take the strictest rule, which in this case is  $x_{U,i} \geq 4$ .

### 3.8 Retrieving a foil data point

Finally, we wish to complement the explanation by presenting the user with a synthetic data point that is similar to the user's data point  $\mathbf{x}_U$ , but that is (correctly) classified as a foil by the foil tree. We refer to such a data point as a *foil data point*. Note that it is possible for samples in a foil leaf to have a classification different from  $B$ , so care needs to be taken in determining the foil sample.

As mentioned in section 3.4, we assume that for each leaf node we saved the secret-shared availability vector  $\xi$  that indicates which data points are present in the leaf node. In section 3.6, we determined the foil leaf, so we can retrieve the corresponding binary availability vector  $\xi^{foil}$ . Recall that the  $i$ -th entry in this vector equals 1, if data point  $\mathbf{x}_i$  is present in the foil leaf, and 0 otherwise. All foil data points  $\mathbf{x}_{i^*}$  therefore satisfy  $\xi_{i^*} = 1$ , and are classified as  $B$  by the foil tree.

A protocol for retrieving a foil data point is presented in Protocol 6. It conceptually works as follows. First, it constructs a indicator vector for the position of the foil data point. This vector is constructed in an element-wise fashion with a helper variable  $\varepsilon$  that indicates whether we already found a foil data point earlier. Second, the secure indicator vector is used to construct the foil data point  $\mathbf{s}$ , which is then revealed to the user.

It is important that the foil data point is only revealed to the user, and not to the computing parties, since the foil data point can leak information on the classifications of the synthetic data points according to the secret-shared model, which are the values we are trying to protect. In practice this means that all

computing parties send their shares of the feature values in vector  $\mathbf{s}$  to the user, who can then combine them to obtain the revealed values.

## 4 Security

We use the MPyC platform [17], which is known to be passively secure. The computing parties jointly form the Explanation Provider (see Figure 1) that securely computes an explanation, which is revealed to the user, who is typically not one of the computing parties. The machine learning model is out of scope, we simply assume secret classifications of synthetic data points are available as secret sharings of the Explanation Provider.

During the protocol, the Explanation Provider will learn the data point  $x_U$  of the user, its class  $A$ , and the foil class  $B$ , together with the average and variance of each feature, used to generate synthetic data set  $\mathcal{X}$ . Furthermore, the (binary) structure of the decision tree, including the fact leaf, foil leaf, and therefore also the lowest common node, will be revealed. Other than this, no training data or model information will be known to the Explanation Provider.

The explanation, consisting of feature index and threshold for each node on the path from lowest common node to fact or foil leaf, and the foil data point  $s$ , is revealed only to the user.

## 5 Complexity

For the generation of the binary decision tree, the number  $\varsigma \approx N \cdot P$  of all possible splits will be large, and determine the runtime. For each node, we need to compute the Gini index for all  $\varsigma$  possibilities, and compute the maximum. If we can compute secure inner products at the cost of one secure multiplication, as in MPyC, the node complexity will be linear in  $\varsigma$  and  $K$ , and more or less equal to the costs of  $\varsigma$  secure comparisons per node. A secure comparison is roughly linear in the number of input bits, which in our case is  $\mathcal{O}(\log_2(NK))$ .

However, we can always precompute  $[y_{i,k}] \cdot [\xi_i]$  for all  $i \in \{1, \dots, N\}$ , and  $k \in \{1, \dots, K\}$ , such that the node complexity is linear in  $N$  and  $K$ . The  $\varsigma$  secure comparisons per node can not be avoided though.

The number of nodes of the decision tree will vary between 1 (no split) and  $\frac{2N-1}{\tau \cdot N}$  (full binary tree). Therefore, the total computational (and communication) complexity will be  $\mathcal{O}(N^2 \cdot K)$ . Although the aim is to obtain a tree of depth  $\log_2 N$ , the depth  $d$  of the tree will vary between 1 (no split) and  $\frac{N-1}{\tau \cdot N}$  (only extremely unbalanced splits). At each tree level we can find the best splits in parallel, such that the number of communication rounds will be limited to  $\mathcal{O}(N \cdot \log_2 \varsigma)$  (assuming a constant round secure comparison).

Given the decision tree, completing the explanation is less complex, and costs at most  $d$  secure comparisons.

$N$	Tree Training			Explanation			Data Point			Accuracy
	avg	min	max	avg	min	max	avg	min	max	
50	20.396	19.594	21.158	0.033	0.027	0.041	0.157	0.112	0.219	0.96
100	94.455	93.133	95.234	0.061	0.058	0.062	0.277	0.269	0.361	0.89
150	130.575	129.681	131.327	0.050	0.038	0.052	0.404	0.387	0.425	0.91

**Table 2.** Performance results (timing in seconds) of our algorithms in MPyC.

## 6 Experiments

We implemented our secure foil tree algorithm in the MPyC framework [17]. This framework functions as a compiler to easily implement protocols that use Shamir secret sharing. It has efficient protocols for scalar-vector multiplications and inner products. In our experiments, we ran MPyC with three parties, and used secure fixed point numbers with a 64-bit integer part and 32-bit fractional part. For the secret-shared black-box model, we secret-shared a neural network with three hidden layers of size 10 each. We used the iris data set [8] as our training data for the neural network (using integer encoding for the target variable) and generated three synthetic data sets based on the iris data set of sizes 50, 100 and 150 respectively.

Table 2 shows the results of our performance tests. We report the timing in seconds of our secure foil tree training algorithm under ‘Tree Training’, for the explanation construction under ‘Explanation’, and for the extraction of the data point under ‘Data Point’. For each of these, we report the average timing, the minimum and the maximum that we observed. The column ‘Accuracy’ denotes the accuracy of the foil tree with respect to the neural network. This accuracy is computed as the number of samples from the synthetic data set for which the classifications according to the neural network and the foil tree are equal, divided by the total number of samples ( $N$ ). We do not provide any performance results on the training algorithm or classification algorithm of the secret-shared black-box model (in this case, the neural network), as the performance of the model highly depends on which model is used, and our solution is model-agnostic.

We see that the accuracy does not necessarily increase when we use more samples. A synthetic data set size of 50 seems to suffice for the iris data set, and shows performance numbers of less than half a minute for the entire algorithm.

## 7 Conclusion

We presented the first cryptographic protocol that is able to explain black-box AI models that are trained by sensitive data, in a privacy-preserving way. The explanation is constructed by means of local foil trees. After generating synthetic data close to a fact data point, a binary tree is securely computed to find the so-called fact and foil leaves. Using both fact and foil leaf, an explanation of the AI model is constructed that explains to the user why they were classified as



the fact class, and not as the foil class. We additionally provide a synthetic data point from the foil leaf to strengthen the explanation.

Our solution hides the classification model and its training data, in order to provide explanations towards users without leaking commercially or privacy sensitive data. We implemented our solution with MPyC on the iris data set with different sizes of synthetic data sets. With 50 samples, we achieved an accuracy of 0.96 within half a minute.

## Acknowledgements

The research of this paper has been done within the FATE project, which is funded by the TNO Appl.AI program (internal AI program). We additionally thank Jasper van der Waa for his helpful comments and suggestions.

## References

1. Mark Abspoel, Daniel Escudero, and Nikolaj Volgushev, *Secure training of decision trees with continuous attributes*, Proc. Priv. Enhancing Technol. **2021** (2021), no. 1, 167–187.
2. Samuel Adams, Chaitali Choudhary, Martine De Cock, Rafael Dowsley, David Melanson, Anderson C. A. Nascimento, Davis Railsback, and Jianwei Shen, *Privacy-preserving training of tree ensembles over continuous data*, CoRR **abs/2106.02769** (2021).
3. Leo Breiman, J. H. Friedman, R. A. Olshen, and C. J. Stone, *Classification and regression trees*, Wadsworth, 1984.
4. Nitesh V. Chawla, Kevin W. Bowyer, Lawrence O. Hall, and W. Philip Kegelmeyer, *SMOTE: synthetic minority over-sampling technique*, J. Artif. Intell. Res. **16** (2002), 321–357.
5. Ronald Cramer, Ivan Damgård, and Jesper Buus Nielsen, *Secure multiparty computation and secret sharing*, Cambridge University Press, 2015.
6. Sebastiaan de Hoogh, *Design of large scale applications of secure multiparty computation: secure linear programming*, Ph.D. thesis, Eindhoven University of Technology, 2012.
7. Sebastiaan de Hoogh, Berry Schoenmakers, Ping Chen, and Harm op den Akker, *Practical secure decision tree learning in a teletreatment application*, Financial Cryptography and Data Security - 18th International Conference, FC 2014, Christ Church, Barbados, March 3-7, 2014, Revised Selected Papers (Nicolas Christin and Reihaneh Safavi-Naini, eds.), Lecture Notes in Computer Science, vol. 8437, Springer, 2014, pp. 179–194.
8. Dheeru Dua and Casey Graff, *UCI machine learning repository*, 2017.
9. Matthew Fredrikson, Eric Lantz, Somesh Jha, Simon M. Lin, David Page, and Thomas Ristenpart, *Privacy in pharmacogenetics: An end-to-end case study of personalized warfarin dosing*, Proceedings of the 23rd USENIX Security Symposium, San Diego, CA, USA, August 20-22, 2014 (Kevin Fu and Jaeyeon Jung, eds.), USENIX Association, 2014, pp. 17–32.
10. Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron C. Courville, and Yoshua Bengio, *Generative adversarial nets*, Advances in Neural Information Processing Systems 27: Annual

- Conference on Neural Information Processing Systems 2014, December 8-13 2014, Montreal, Quebec, Canada (Zoubin Ghahramani, Max Welling, Corinna Cortes, Neil D. Lawrence, and Kilian Q. Weinberger, eds.), 2014, pp. 2672–2680.
11. Frederik Harder, Matthias Bauer, and Mijung Park, *Interpretable and differentially private predictions*, The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020, AAAI Press, 2020, pp. 4083–4090.
  12. Bart Kamphorst, Daan Knoors, and Thomas Rooijackers, *Oncological research on distributed patient data: Privacy can be preserved!*, ERCIM News **2021** (2021), no. 126.
  13. Scott M. Lundberg and Su-In Lee, *A unified approach to interpreting model predictions*, Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA (Isabelle Guyon, Ulrike von Luxburg, Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, and Roman Garnett, eds.), 2017, pp. 4765–4774.
  14. Pascal Paillier, *Public-key cryptosystems based on composite degree residuosity classes*, Advances in Cryptology - EUROCRYPT '99, International Conference on the Theory and Application of Cryptographic Techniques, Prague, Czech Republic, May 2-6, 1999, Proceeding (Jacques Stern, ed.), Lecture Notes in Computer Science, vol. 1592, Springer, 1999, pp. 223–238.
  15. Marco Túlio Ribeiro, Sameer Singh, and Carlos Guestrin, *"why should I trust you?": Explaining the predictions of any classifier*, Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, San Francisco, CA, USA, August 13-17, 2016 (Balaji Krishnapuram, Mohak Shah, Alexander J. Smola, Charu C. Aggarwal, Dou Shen, and Rajeev Rastogi, eds.), ACM, 2016, pp. 1135–1144.
  16. Alex Sangers, Maran van Heesch, Thomas Attema, Thijs Veugen, Mark Wiggerman, Jan Veldsink, Oscar Bloemen, and Daniël Worm, *Secure multiparty pagerank algorithm for collaborative fraud detection*, Financial Cryptography and Data Security (Cham) (Ian Goldberg and Tyler Moore, eds.), Springer International Publishing, 2019, pp. 605–623.
  17. Berry Schoenmakers, *MPC - Secure Multiparty Computation in Python*, <https://github.com/lschoe/mpyc>.
  18. Tomas Toft, *Solving linear programs using multiparty computation*, Financial Cryptography and Data Security, 13th International Conference, FC 2009, Accra Beach, Barbados, February 23-26, 2009. Revised Selected Papers (Roger Dingledine and Philippe Golle, eds.), Lecture Notes in Computer Science, vol. 5628, Springer, 2009, pp. 90–107.
  19. Jasper van der Waa, Elisabeth Nieuwburg, Anita H. M. Cremers, and Mark A. Neerincx, *Evaluating XAI: A comparison of rule-based and example-based explanations*, Artif. Intell. **291** (2021), 103404.
  20. Jasper van der Waa, Marcel Robeer, Jurriaan van Diggelen, Matthieu Brinkhuis, and Mark A. Neerincx, *Contrastive explanations with local foil trees*, CoRR **abs/1806.07470** (2018).
  21. Marie Beth van Egmond, Gabriele Spini, Onno van der Galiën, Arne IJpma, Thijs Veugen, Wessel Kraaij, Alex Sangers, Thomas Rooijackers, Peter Langenkamp, Bart Kamphorst, Natasja van de L'Isle, and Milena Kooij-Janic, *Privacy-preserving*

- dataset combination and lasso regression for healthcare predictions*, BMC Medical Informatics Decis. Mak. **21** (2021), no. 1, 266.
22. Thijs Veugen, Bart Kamphorst, Natasja van de L’Isle, and Marie Beth van Egmond, *Privacy-preserving coupling of vertically-partitioned databases and subsequent training with gradient descent*, Cyber Security Cryptography and Machine Learning - 5th International Symposium, CSCML 2021, Be’er Sheva, Israel, July 8-9, 2021, Proceedings (Shlomi Dolev, Oded Margalit, Benny Pinkas, and Alexander A. Schwarzmann, eds.), Lecture Notes in Computer Science, vol. 12716, Springer, 2021, pp. 38–51.
  23. Ziqi Yang, Jiyi Zhang, Ee-Chien Chang, and Zhenkai Liang, *Neural network inversion in adversarial setting via background knowledge alignment*, Proceedings of the 2019 ACM SIGSAC Conference on Computer and Communications Security, CCS 2019, London, UK, November 11-15, 2019 (Lorenzo Cavallaro, Johannes Kinder, XiaoFeng Wang, and Jonathan Katz, eds.), ACM, 2019, pp. 225–240.
  24. Andrew C. Yao, *Protocols for secure computations*, 23rd Annual Symposium on Foundations of Computer Science (sfcs 1982), 1982, pp. 160–164.
  25. Yuheng Zhang, Ruoxi Jia, Hengzhi Pei, Wenxiao Wang, Bo Li, and Dawn Song, *The secret revealer: Generative model-inversion attacks against deep neural networks*, 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2020, Seattle, WA, USA, June 13-19, 2020, Computer Vision Foundation / IEEE, 2020, pp. 250–258.

## A Indicator vector of maximum

Given a list of hidden elements  $[z] = ([z_1], \dots, [z_L])$  and a relation  $<$  that induces a total ordering on its elements, we need to find the maximum  $[z_{\max}]$  and the indicator vector of the maximum. One way to do this is to securely deduce  $[z_{\max}]$  and then apply an independent protocol for finding the position of  $[z_{\max}]$  and returning the result as an indicator vector. This is supported out-of-the-box in several frameworks.

Instead, we suggest to store some artefacts of the first protocol and leverage them in the second protocol. This is achieved through Protocols 7 and 8. First,  $\max([z], 1, L)$  performs a binary search to compute the maximum through  $L - 1$  secure comparisons, which comparison results are stored in the  $[\gamma_s]$ ,  $1 \leq s < L$ , followed by  $\text{indicator}([1], 1, L)$  to compute the indicators  $[\delta_s]$ ,  $1 \leq s \leq L$  of the maximum.

This approach with logarithmic round complexity is similar to Protocol 5.1 of de Hoogh [6], and due to Toft [18]. Since the both recursive calls in line 7.5 can be performed in parallel, the number of iterations is reduced from  $L$  to  $\log_2 L$ .

## B Comparing adjusted Gini indices

Given two splits  $S_1$  and  $S_2$ , we wish to compare their adjusted Gini indices  $[\hat{G}_1] = [\tilde{N}_1]/[\tilde{D}_1]$  and  $[\hat{G}_2] = [\tilde{N}_2]/[\tilde{D}_2]$ . In particular, we wish to compute  $[(\hat{G}_1 < \hat{G}_2)]$  with the interpretation that an index with zero-valued denominator is always smaller than the index. If both denominators are zero, the result does not matter.

**Protocol 7 max**

Computes the maximum

**Input:** Vector  $[z]$ , indices  $s_l$  and  $s_r$ **Output:** Maximum  $[\max\{z_s \mid s_l \leq s \leq s_r\}]$ , storing comparison results  $[\gamma_s]$ ,  
 $s_l \leq s < s_r$ 

- 1: **if**  $s_l = s_r$  **then**
- 2:      $[z_{\max}] \leftarrow [z_{s_l}]$
- 3: **else**
- 4:      $\bar{s} \leftarrow (s_l + s_r) \div 2$  ▷ split at  $\bar{s}$ ,  $s_l \leq \bar{s} < s_r$
- 5:      $[z_l] \leftarrow \max([z], s_l, \bar{s})$ ;  $[z_r] \leftarrow \max([z], \bar{s} + 1, s_r)$
- 6:      $[\gamma_{\bar{s}}] \leftarrow [(z_l < z_r)]$  ▷ Is  $z_r$  the largest?
- 7:      $[z_{\max}] \leftarrow [z_l] + [\gamma_{\bar{s}}] \cdot ([z_r] - [z_l])$
- 8: **end if**
- 9: Return  $[z_{\max}]$

**Protocol 8 indicator**

Secure maximum indicator vector

**Input:** Indicator  $[\delta]$ , indices  $s_l$  and  $s_r$ , and comparison results  $[\gamma]$ **Output:** Indicators  $\{[\delta_s] \mid s_l \leq s \leq s_r\}$  of the overall maximum**Invariant:**  $\delta = (s_l \leq \operatorname{argmax}\{z_s \mid 1 \leq s \leq L\} \leq s_r)$ 

- 1: **if**  $s_l = s_r$  **then**
- 2:     Return  $[\delta]$
- 3: **else**
- 4:      $\bar{s} \leftarrow (s_l + s_r) \div 2$  ▷ split at  $\bar{s}$ ,  $s_l \leq \bar{s} < s_r$
- 5:      $[\epsilon] \leftarrow [\delta] \cdot [\gamma_{\bar{s}}]$
- 6:      $[\delta_l] \leftarrow \operatorname{indicator}([\delta] - [\epsilon], s_l, \bar{s})$ ;  $[\delta_r] \leftarrow \operatorname{indicator}([\epsilon], \bar{s} + 1, s_r)$
- 7:     Return the elements of  $[\delta_l]$  and  $[\delta_r]$
- 8: **end if**

To avoid complications when either denominator is zero, we change the straightforward integer comparison  $N_1 \cdot D_2 < N_2 \cdot D_1$  to

$$N_1 \cdot (N_1 \cdot D_2 - D_1 \cdot N_2) < 1 - D_1. \quad (5)$$

To see why this is correct, recall that  $N_s$  and  $D_s$  are both integer, and  $D_s = 0$ , if and only if,  $N_s = 0$ . Additionally, it follows from Equation (4) that

$$N_s \geq n^r \sum_{k=1}^K n_k^l + n^l \sum_{k=1}^K n_k^r = 2D_s. \quad (6)$$

Therefore, the following statements hold:

- If  $D_1, D_2 > 0$ , then  $N_1, N_2 > 0$  and (5) evaluates to the result of  $N_1 \cdot D_2 - D_1 \cdot N_2 < \frac{1-D_1}{N_1} \in (-1/2, 0]$ . Since all variables on the left-hand side are integers, this is equivalent to  $N_1 \cdot D_2 - D_1 \cdot N_2 < 0$ .
- If  $D_1 > 0, D_2 = 0$ , then  $N_2 = 0$  and (5) evaluates to the result of  $0 < 1 - D_1$ , which is False.

- If  $D_1 = 0$ , then  $N_1 = 0$  and (5) evaluates to the result of  $0 < 1$ , which is True.

An alternative approach is to compute  $N_1 \cdot (D_2 + 1) < N_2 \cdot (D_1 + 1)$ . Theoretically this comparison might<sup>5</sup> not indicate the worst adjusted Gini index if the indices have very small difference, but a significant efficiency boost can be expected as the secure comparison input can be represented in fewer bits.

---

<sup>5</sup> We did not prove this, the opposite could be proven if some convenient properties are derived from Equation (4).