



## **SIRUS: Stable and Interpretable RULE Set for Classification**

Clément Bénard, Gérard Biau, Sébastien Da Veiga, Erwan Scornet

### **► To cite this version:**

Clément Bénard, Gérard Biau, Sébastien Da Veiga, Erwan Scornet. SIRUS: Stable and Interpretable RULE Set for Classification. *Electronic Journal of Statistics*, 2021, 15 (1), pp.427 - 505. <10.1214/20-EJS1792>. <hal-02190689v4>

**HAL Id: hal-02190689**

**<https://hal.science/hal-02190689v4>**

Submitted on 15 Dec 2020

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



HAL Authorization

# SIRUS: Stable and Interpretable RULe Set for Classification

Clément Bénard\*   Gérard Biau<sup>†</sup>   Sébastien Da Veiga<sup>‡</sup>   Erwan Scornet<sup>§</sup>

## Abstract

State-of-the-art learning algorithms, such as random forests or neural networks, are often qualified as “black-boxes” because of the high number and complexity of operations involved in their prediction mechanism. This lack of interpretability is a strong limitation for applications involving critical decisions, typically the analysis of production processes in the manufacturing industry. In such critical contexts, models have to be interpretable, i.e., simple, stable, and predictive. To address this issue, we design SIRUS (Stable and Interpretable RULe Set), a new classification algorithm based on random forests, which takes the form of a short list of rules. While simple models are usually unstable with respect to data perturbation, SIRUS achieves a remarkable stability improvement over cutting-edge methods. Furthermore, SIRUS inherits a predictive accuracy close to random forests, combined with the simplicity of decision trees. These properties are assessed both from a theoretical and empirical point of view, through extensive numerical experiments based on our R/C++ software implementation `sirus` available from `CRAN`.

**Keywords:** classification, interpretability, rules, stability, random forests.

## 1 Introduction

State-of-the-art learning algorithms, typically tree ensembles or neural networks, are well-known for their remarkable predictive performance. However, this high accuracy comes at the price of complex prediction mechanisms: a large number of operations are computed for a given prediction. Because of this complexity, learning algorithms are often considered as black-boxes. This lack of interpretability is a serious limitation for many applications involving critical decisions, such as healthcare, criminal justice, or industrial process optimization. This latter example is interesting to illustrate how interpretability can be essential. Indeed, in the manufacturing industry, production processes involve complex physical and chemical phenomena, whose control and efficiency are of critical importance. In practice, data is collected along the manufacturing line, describing both the production environment and its conformity. The retrieved information enables to infer a link between the manufacturing conditions and

---

\*Safran Tech, Sorbonne Université

<sup>†</sup>Sorbonne Université

<sup>‡</sup>Safran Tech

<sup>§</sup>Ecole Polytechnique

the resulting quality at the end of the line, and then to increase the process efficiency. Since the quality of the produced entities is often characterized by a pass or fail output, the problem is in fact a classification task, and state-of-the-art learning algorithms can successfully catch patterns of these complex and nonlinear physical phenomena. However, any decision impacting the production process has long-term and heavy consequences, and therefore cannot simply rely on a blind stochastic modelling. As a matter of fact, a deep physical understanding of the forces in action is required, and this makes black-box algorithms inappropriate. In a word, models have to be interpretable, i.e., provide an understanding of the internal mechanisms that build a relation between inputs and outputs, to provide insights to guide the physical analysis. This is for example typically the case in the aeronautics industry, where the manufacturing of engine parts involves sensitive casting and forging processes. Interpretable models allow us to gain knowledge on the behavior of such production processes, which can lead, for instance, to identify or fine-tune critical parameters, improve measurement and control, optimize maintenance, or deepen understanding of physical phenomena. In the following paragraphs, we deepen the discussion about the definition of interpretability to highlight the limitations of the most popular interpretable nonlinear models: decision trees and rule algorithms (Guidotti et al., 2018). Despite their high predictivity and simple structure, these methods are unstable, which is a strong operational limitation. The goal of this article is to introduce **SIRUS** (Stable and Interpretable **R**Ule **S**et), an interpretable rule classification algorithm which considerably improves stability over state-of-the-art methods, while preserving their simple structure, accuracy, and computational complexity.

As stated in Rüping (2006), Lipton (2016), Doshi-Velez and Kim (2017), or Murdoch et al. (2019), to date, there is no agreement in statistics and machine learning communities about a rigorous definition of interpretability. There are multiple concepts behind it, many different types of methods, and a strong dependence on the area of application and the audience. Here, we focus on models intrinsically interpretable, which directly provide insights on how inputs and outputs are related, as opposed to the post-processing of black-box models. In that case, we argue that it is possible to define minimum requirements for interpretability through the triptych “simplicity, stability, and predictivity”, in line with the framework recently proposed by Yu and Kumbier (2019). Indeed, in order to grasp how inputs and outputs are related, the structure of the model has to be simple. The notion of simplicity is implied whenever interpretability is invoked (e.g., Rüping, 2006; Freitas, 2014; Letham, 2015; Letham et al., 2015; Lipton, 2016; Ribeiro et al., 2016; Murdoch et al., 2019) and essentially refers to the model size, complexity, or the number of operations performed in the prediction mechanism. Yu (2013) defines stability as another fundamental requirement for interpretability: conclusions of a statistical analysis have to be robust to small data perturbations to be meaningful. Indeed, a specific analysis is likely to be run multiple times, eventually adding a small new batch of data, and an interpretable algorithm should be insensitive to such modifications. Otherwise, unstable models provide us with a partial and arbitrary analysis of the underlying phenomena, and arouses distrust of the domain experts. Finally, if the predictive accuracy of an interpretable model is significantly lower than the one of a state-of-the-art black-box algorithm, it clearly misses strong patterns in the data and will therefore be useless, as explained in Breiman (2001b). For example, the trivial model that outputs the empirical mean of the observations for any input is simple, stable, but brings in most cases no useful information. Thus, we add a good predictivity as an essential requirement for interpretability.

Decision trees are a class of supervised learning algorithms that recursively partition the input space and make local decisions in the cells of the resulting partition. Trees can model highly nonlinear patterns while having a simple structure, and are therefore good candidates when interpretability is required. However, trees are unstable to small data perturbations (Oates and Jensen, 1997; Guidotti and Ruggieri, 2019). More precisely, as explained in Breiman (2001b): by randomly removing only 2 – 3% of the training data, the tree structure can be quite different, which is a strong limitation to their practical use. Another class of supervised learning methods that can model nonlinear patterns while retaining a simple structure are the so-called rule models. As such, a rule is defined as a conjunction of constraints on input variables, which form a hyperrectangle in the input space where the estimated output is constant. A collection of rules is combined to form a model. Here, the term “rule” does not stand for “classification rule” but, as is traditional in the rule learning literature, to a piecewise constant estimate that simply reads “if *conditions on  $\mathbf{x}$* , then *response*, else *default response*”. Despite their simplicity and excellent predictive skills, rule algorithms are unstable and, from this point of view, share the same limitation as decision trees (Letham et al., 2015; Murdoch et al., 2019).

In line with the above, we design SIRUS in the present paper, a new rule classification algorithm which inherits an accuracy close to random forests and the simplicity of decision trees, while having a stable structure. The core aggregation principle of random forests is kept, but instead of aggregating predictions, SIRUS focuses on the probability that a given hyperrectangle (i.e., a node) is contained in a randomized tree. The nodes with the highest probability are robust to data perturbation and represent strong patterns. They are therefore selected to form a stable rule ensemble model. Here, we provide a first illustration of SIRUS with a simple and real case: the Titanic dataset (Piech, 2016). The survival status of 887 passengers are recorded, as well as various personal characteristics: age, sex, class, number of siblings and parents aboard, and the paid fare. SIRUS outputs the following simple set of 7 rules, which enables to grasp at a glance the main patterns to explain passenger survival:

Average survival rate $p_s = 39\%$ .				
if	sex is male	then	$p_s = 19\%$	else $p_s = 74\%$
if	1 <sup>st</sup> or 2 <sup>nd</sup> class	then	$p_s = 56\%$	else $p_s = 24\%$
if	1 <sup>st</sup> or 2 <sup>nd</sup> class & sex is female	then	$p_s = 95\%$	else $p_s = 25\%$
if	fare < 10.5£	then	$p_s = 20\%$	else $p_s = 50\%$
if	no parents or children aboard	then	$p_s = 35\%$	else $p_s = 51\%$
if	2 <sup>st</sup> or 3 <sup>rd</sup> class & sex is male	then	$p_s = 14\%$	else $p_s = 64\%$
if	sex is male & age $\geq 15$	then	$p_s = 16\%$	else $p_s = 72\%$

To generate the prediction for a new query point  $\mathbf{x}$ , SIRUS checks for each rule whether the conditions are satisfied to assign one of the two possible  $p_s$  output values. Let us say for

example that  $x^{(sex)}$  is female, then  $\mathbf{x}$  satisfies the condition of the first rule, which returns  $p_s = 74\%$ . Next, the 7 rule outputs are averaged to provide the predicted probability of survival for  $\mathbf{x}$ . The model is stable: when a 10-fold cross-validation is run to simulate data perturbation, 5 to 6 rules are consistent across two folds in average. The model error (1-AUC) is 0.17, close to the 0.13 of random forests, whereas simplicity is drastically increased: 7 rules versus about  $10^4$  operations for a forest prediction.

First, we review the main rule algorithms and present their mechanism principles in Section 2. Next, Section 3 is devoted to the detailed description of SIRUS. One of the main contributions of this work is the development of a software implementation, via the R/C++ package `sirus` (Benard and Wright, 2020) available from CRAN, based on `ranger`, a high-performance random forest implementation (Wright and Ziegler, 2017). In Section 4, we show that the good empirical behavior of SIRUS is theoretically understood by proving its asymptotic stability. Then, in Section 5, we illustrate the efficiency of our algorithm through numerical experiments on real datasets. Finally, Section 6 summarizes the main contributions of the article and provides directions for future research.

## 2 Related Work

As stated in the introduction, SIRUS has two types of competitors: decision trees and rule algorithms. More precisely, the latter can further be split into three different kinds: classical rule algorithms based on greedy heuristics, those built on top of frequent pattern mining algorithms, and those extracted from tree ensembles.

Decision trees may be the most popular competitors of SIRUS because of their simple structure. The main algorithms are CART (Breiman et al., 1984) and C5.0 (Quinlan, 1992). However, trees are unstable as we have already highlighted. A widespread method to stabilize decision trees is bagging (Breiman, 1996), in which multiple trees are grown on perturbed data and aggregated together. Random forests is an algorithm developed by Breiman (2001a) that improves over bagging by randomizing the tree construction. Predictions are stable, accuracy is increased, but the final model is unfortunately a black box. Thus, simplicity of trees is lost, and some post-treatment mechanisms are needed to understand how random forests make their decisions. Nonetheless, even if they are useful, such treatments only provide partial information and can be difficult to operationalize for critical decisions (Rudin, 2018). For example, variable importance (Breiman, 2001a, 2003a) identifies variables that have a strong impact on the output, but not which inputs values are associated to output values of interest. Similarly, local approximation methods such as LIME (Ribeiro et al., 2016) or Tolomei et al. (2017) do not provide insights on the global relation.

Rule learning originates from the influential AQ system of Michalski (1969). Many algorithms based on greedy heuristics were subsequently developed in the 1980’s and 1990’s, including Decision List (Rivest, 1987), CN2 (Clark and Niblett, 1989), FOIL (First-Order Inductive Learner, Quinlan, 1990; Quinlan and Cameron-Jones, 1995), IREP (Incremental Reduced Error Pruning, Fürnkranz and Widmer, 1994), RIPPER (Repeated Incremental Pruning to Produce Error Reduction, Cohen, 1995), PART (Partial Decision Trees, Frank and Witten, 1998), SLIPPER (Simple Learner with Iterative Pruning to Produce Error Reduction, Cohen

and Singer, 1999), LRI (Leightweight Rule Induction, Weiss and Indurkha, 2000), and EN-DER (Ensemble of Decision Rules, Dembczyński et al., 2010). Since these methods are based on greedy heuristics, they are computationally fast, but similarly to decision trees, they are unstable and their accuracy is often limited.

At the end of the 1990’s a new type of rule algorithms based on frequent pattern mining is introduced with CBA (Classification Based on Association Rules, Liu et al., 1998), then extended with CPAR (Classification based on Predictive Association Rules, Yin and Han, 2003). Frequent pattern mining is originally used to identify frequent occurrences in database mining. Since the output  $Y \in \{0, 1\}$  is discrete and the input data can be discretized, we can generate candidate rules for classification by identifying frequent patterns associated with each output label. This exhaustive search for association rules is computationally costly (exponential with the input dimension), and efficient heuristics are used, essentially Apriori (Agrawal et al., 1993) and Eclat (Zaki et al., 1997). The rule aggregation mechanism is specific to each algorithm. More recently, BRL (Bayesian Rule List, Letham et al., 2015) uses a more sophisticated Bayesian framework for the rule aggregation than the simple approach of CBA and CPAR, while IDS (Lakkaraju et al., 2016, Interpretable Decision Sets) uses a multi-objective optimization to select interpretable rules. Finally, CORELS (Angelino et al., 2017, Certifiably Optimal Rule ListS) generates optimal rule lists for categorical data. Interestingly, these methods exhibit quite good stability properties as we will see, higher than decision trees, but on the other hand, their predictive accuracy is worse.

The last decade has seen a resurgence of rule models through powerful algorithms based on rule extraction from tree ensembles, especially RuleFit (Friedman and Popescu, 2008) and Node harvest (Meinshausen, 2010). Notice that SIRUS is also based on this principle. More specifically, RuleFit extracts all the rules of a boosted tree ensemble (Friedman and Popescu, 2003), while Node harvest is based on random forests. Then, the extracted rules are linearly combined in a sparse linear model, respectively a logistic regression with a Lasso penalty (Tibshirani, 1996) for RuleFit, and a constraint quadratic linear program for Node harvest. These two methods have a computational complexity comparable to random forests and SIRUS, since the main step of all these algorithms is to grow a tree ensemble with a large number of trees. However, both algorithms are unstable, and both output quite complex and long lists of rules. Even running RuleFit or Node harvest multiple times on the same dataset produces quite different rule lists because of the randomness in the tree ensembles—see Appendix A.1. On the other hand, SIRUS is built to have its structure converged for the given dataset, as explained later in Section 3.

To the best of our knowledge, the signed iterative random forest method (s-iRF, Kumbier et al., 2018) is the only procedure that tackles both rule learning and stability. Using random forests, s-iRF manages to extract stable signed interactions, i.e., feature interactions enriched with a thresholding behavior for each variable, lower or higher, but without specific thresholding values. Therefore, s-iRF can be difficult to operationalize since it does not provide any specific input thresholds, and thus no precise information about the influence of input variables. On the other hand, an explicit rule model identifies specific regions of interest in the input space.

### 3 SIRUS Algorithm

Within the general framework of supervised (binary) classification, we assume to be given an i.i.d. sample  $\mathcal{D}_n = \{(\mathbf{X}_i, Y_i), i = 1, \dots, n\}$ . Each  $(\mathbf{X}_i, Y_i)$  is distributed as the generic pair  $(\mathbf{X}, Y)$  independent of  $\mathcal{D}_n$ , where  $\mathbf{X} = (X^{(1)}, \dots, X^{(p)})$  is a random vector taking values in  $\mathbb{R}^p$  and  $Y \in \{0, 1\}$  is a binary response. Throughout the document, the distribution of  $(\mathbf{X}, Y)$  is assumed to be unknown and is denoted by  $\mathbb{P}_{\mathbf{X}, Y}$ . For  $\mathbf{x} \in \mathbb{R}^p$ , our goal is to accurately estimate the conditional probability  $\eta(\mathbf{x}) = \mathbb{P}(Y = 1 | \mathbf{X} = \mathbf{x})$  with few simple and stable rules.

To tackle this problem, SIRUS first builds a (slightly modified) random forest. Next, each hyperrectangle of each tree of the forest is turned into a simple decision rule, and the collection of these elementary rules is ranked based on their frequency of appearance in the forest. Finally, the most significant rules are retained and are averaged together to form an ensemble model. We describe the four steps of SIRUS algorithm in the following paragraphs: the rule generation, rule selection, rule post-treatment, and the rule aggregation. This section ends with a discussion of SIRUS stability.

**Rule generation.** SIRUS uses at its core the random forest method (Breiman, 2001a), slightly modified for our purpose. As in the original procedure, each single tree in the forest is grown with a greedy heuristic that recursively partitions the input space using a random variable  $\Theta$ . The essential difference between our approach and Breiman’s one is that, prior to all tree constructions, the empirical  $q$ -quantiles of the marginal distributions over the whole dataset are computed: in each node of each tree, the best split can be selected among these empirical quantiles only. This constraint is critical to stabilize the forest structure and keeps almost intact the predictive accuracy, provided  $q$  is not too small (typically of the order of 10—see the experimental Subsection 5.4). Apart from this difference, the tree growing is similar to Breiman’s original procedure. The tree randomization  $\Theta$  is independent of the sample and has two independent components, denoted by  $\Theta^{(S)}$  and  $\Theta^{(V)}$ , which are respectively used for the subsampling mechanism and randomization of the split direction. Throughout the manuscript, we let  $\hat{q}_{n,r}^{(j)}$  be the empirical  $r$ -th  $q$ -quantile of  $\{X_1^{(j)}, \dots, X_n^{(j)}\}$ , with typically  $q = 10$ . The construction of the individual trees is summarized in Algorithm 1 below.

---

#### Algorithm 1 Tree construction

---

- 1: **Parameters:** Number of quantiles  $q$ , number of subsampled observations  $a_n$ , number of eligible directions for splitting `mtry`.
  - 2: Compute the empirical  $q$ -quantiles for each marginal distribution over the whole dataset.
  - 3: Subsample with replacement  $a_n$  observations, indexed by  $\Theta^{(S)}$ . Only these observations are used to build the tree.
  - 4: Initialize the cell  $H$  as the root of the tree.
  - 5: Draw uniformly at random a subset  $\Theta^{(V)} \subset \{1, \dots, p\}$  of cardinality `mtry`.
  - 6: For all  $j \in \Theta^{(V)}$ , compute the CART-splitting criterion at all empirical  $q$ -quantiles of  $X^{(j)}$  that split the cell  $H$  into two non-empty cells.
  - 7: Choose the split that maximizes the CART-splitting criterion.
  - 8: Recursively repeat **lines** 5 – 7 for the two resulting children cells  $H_L$  and  $H_R$ .
-



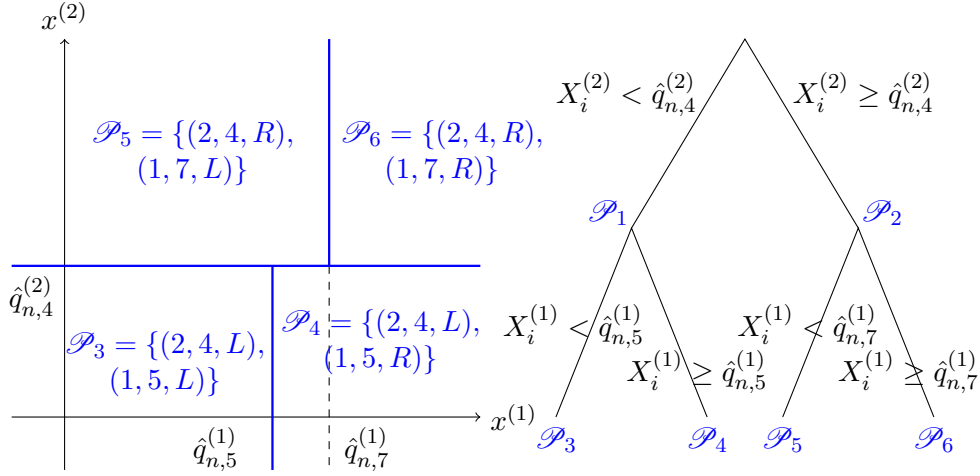


Figure 1: Example of a root node  $\mathbb{R}^2$  partitionned by a randomized tree of depth 2: the tree on the right side, the associated paths and hyperrectangles of length  $d = 2$  on the left side.

The main step of SIRUS is to extract rules from the modified random forest. The cornerstone of this extraction mechanism is the notion of path in a decision tree. Indeed, a path describes the sequence of splits to go from the root of the tree to a specific (inner or terminal) node. Since a hyperrectangle is associated to each node, a rule can be defined as a piecewise constant estimate with this hyperrectangle as support. Therefore, to rigorously define the rule extraction, we introduce the symbolic representation of a path in a tree. We insist that such definition is valid for both terminal leaves and inner nodes, which are all used by SIRUS. To begin, we follow the example shown in Figure 1 with a tree of depth 2 partitioning the input space  $\mathbb{R}^2$ . For instance, let us consider the node  $\mathcal{P}_6$  defined by the sequence of two splits  $X_i^{(2)} \geq \hat{q}_{n,4}^{(2)}$  and  $X_i^{(1)} \geq \hat{q}_{n,7}^{(1)}$ . The first split is symbolized by the triplet  $(2, 4, R)$ , whose components respectively stand for the variable index 2, the quantile index 4, and the right side  $R$  of the split. Similarly, for the second split we cut coordinate 1 at quantile index 7, and pass to the right. Thus, the path to the considered node is defined by  $\mathcal{P}_6 = \{(2, 4, R), (1, 7, R)\}$ . Also notice that the first split already defines the path  $\mathcal{P}_2 = \{(2, 4, R)\}$ , associated to the right inner node at the first level of the tree. Of course, this generalizes to each path  $\mathcal{P}$  of length  $d$  under the symbolic compact form

$$\mathcal{P} = \{(j_k, r_k, s_k), k = 1, \dots, d\},$$

where, for  $k \in \{1, \dots, d\}$ , the triplet  $(j_k, r_k, s_k)$  describes how to move from level  $(k - 1)$  to level  $k$ , with a split using the coordinate  $j_k \in \{1, \dots, p\}$ , the index  $r_k \in \{1, \dots, q - 1\}$  of the corresponding quantile, and a side  $s_k = L$  if we go to the left and  $s_k = R$  if we go to the right. The set of all possible such paths is denoted by  $\Pi$ . It is important to note that  $\Pi$  is in fact a deterministic (that is, non random) quantity, which only depends upon the dimension  $p$  and the order  $q$  of the quantiles. Of course, given a path  $\mathcal{P} \in \Pi$  one can recover the hyperrectangle (i.e., the tree node)  $\hat{H}_n(\mathcal{P})$  associated with  $\mathcal{P}$  and the entire dataset  $\mathcal{D}_n$



via the correspondence

$$\hat{H}_n(\mathcal{P}) = \left\{ \mathbf{x} \in \mathbb{R}^p : \begin{cases} \mathbf{x}^{(j_k)} < \hat{q}_{n,r_k}^{(j_k)} & \text{if } s_k = L \\ \mathbf{x}^{(j_k)} \geq \hat{q}_{n,r_k}^{(j_k)} & \text{if } s_k = R \end{cases}, k = 1, \dots, d \right\}. \quad (3.1)$$

Finally, an elementary rule  $\hat{g}_{n,\mathcal{P}}$  can be defined from  $\hat{H}_n(\mathcal{P})$  as a piecewise constant estimate:  $\hat{g}_{n,\mathcal{P}}(\mathbf{x})$  returns the empirical probability that the output  $Y$  is of class 1 conditional on whether the query point  $\mathbf{x}$  belongs to  $\hat{H}_n(\mathcal{P})$  or not. Thus, the rule  $\hat{g}_{n,\mathcal{P}}$  associated to the path  $\mathcal{P} \in \Pi$  is formally defined by

$$\forall \mathbf{x} \in \mathbb{R}^p, \quad \hat{g}_{n,\mathcal{P}}(\mathbf{x}) = \begin{cases} \frac{1}{N_n(\hat{H}_n(\mathcal{P}))} \sum_{i=1}^n Y_i \mathbb{1}_{\mathbf{x}_i \in \hat{H}_n(\mathcal{P})} & \text{if } \mathbf{x} \in \hat{H}_n(\mathcal{P}) \\ \frac{1}{n - N_n(\hat{H}_n(\mathcal{P}))} \sum_{i=1}^n Y_i \mathbb{1}_{\mathbf{x}_i \notin \hat{H}_n(\mathcal{P})} & \text{otherwise} \end{cases},$$

using the convention  $0/0 = 0$ , and where  $N_n(\hat{H}_n(\mathcal{P}))$  is the number of observations in the node associated with  $\mathcal{P}$ . This formal definition can be illustrated with the Titanic dataset presented in the introduction. For the fourth rule, `fare` is the 6th variable and since  $\hat{q}_{n,4}^{(6)} = 10.5$ , the corresponding path is  $\mathcal{P} = \{(6, 4, L)\}$ , and the associated rule is thus

$$\hat{g}_{n,\mathcal{P}}(\mathbf{x}) = \begin{cases} 0.20 & \text{if } x^{(6)} < 10.5 \\ 0.50 & \text{if } x^{(6)} \geq 10.5 \end{cases}.$$

Finally, a  $\Theta$ -random tree generates a collection of paths in  $\Pi$ , one for each internal and terminal nodes. In the sequel, we let  $T(\Theta, \mathcal{D}_n)$  be the list of such extracted paths, a random subset of  $\Pi$ .

**Rule selection.** Using our modified random forest algorithm, we are able to generate a large number  $M$  of trees, randomized by  $\Theta_1, \dots, \Theta_M$ , i.i.d. copies of the generic variable  $\Theta$ , and then to extract a large collection of rules. Since we are interested in selecting the most important rules, i.e., those which represent strong patterns between the inputs and the output, we select rules that are shared by a large portion of trees. Such occurrence frequency is formally defined by

$$\hat{p}_{M,n}(\mathcal{P}) = \frac{1}{M} \sum_{\ell=1}^M \mathbb{1}_{\mathcal{P} \in T(\Theta_\ell, \mathcal{D}_n)},$$

which is the Monte-Carlo estimate of the probability that a path  $\mathcal{P}$  belongs to a  $\Theta$ -random tree, that is

$$p_n(\mathcal{P}) = \mathbb{P}(\mathcal{P} \in T(\Theta, \mathcal{D}_n) | \mathcal{D}_n).$$

As a general strategy, once the modified random forest has been built, we draw the list of all paths that appear in the forest and only retain those that occur with a frequency larger than the threshold  $p_0 \in (0, 1)$ , the only influential parameter of SIRUS—see Subsection 5.4 for its tuning procedure. We are thus interested in the set of the extracted paths

$$\hat{\mathcal{P}}_{M,n,p_0} = \{\mathcal{P} \in \Pi : \hat{p}_{M,n}(\mathcal{P}) > p_0\}. \quad (3.2)$$

An important feature of SIRUS algorithm is to stop the growing of the forest with an appropriate number of trees  $M$ . Although the right order of magnitude for  $M$  is required, no fine

tuning is necessary. Indeed, the uncertainty of the importance estimate  $\hat{p}_{M,n}(\mathcal{P})$  of each rule decreases with  $M$ , whereas the computational cost linearly increases with  $M$ . Thus, to obtain a robust rule extraction,  $M$  needs to be high enough to make the uncertainty of  $\hat{p}_{M,n}(\mathcal{P})$  negligible. More precisely,  $M$  is set to get the same list of selected rules  $\hat{\mathcal{P}}_{M,n,p_0}$  when SIRUS is run multiple times on the same dataset  $\mathcal{D}_n$ . On the other hand,  $M$  should be small enough to avoid useless computations. Therefore, the growing of the forest is automatically stopped when 95% of the selected rules would be shared by a new run of SIRUS on  $\mathcal{D}_n$  in average, as it is possible to derive a simple stopping criterion based on the properties of the estimates  $\hat{p}_{M,n}(\mathcal{P})$ —all the technical details are provided in Subsection 5.4. A random forest is usually built with around 500 trees, as the predictive accuracy cannot be significantly increased by adding more trees. SIRUS typically grows 10 times more trees to obtain a robust rule extraction.

Besides, we insist that the quantile discretization is critical for the rule selection. The expected value of the rule importance is

$$\mathbb{E}[\hat{p}_{M,n}(\mathcal{P})] = \mathbb{P}(\mathcal{P} \in T(\Theta, \mathcal{D}_n)),$$

but without the discretization, the list of extracted paths from a random tree  $T(\Theta, \mathcal{D}_n)$  takes values in an uncountable space when at least one component of  $\mathbf{X}$  is a continuous random variable, and therefore the above quantity is null, making the path selection procedure unstable with respect to data perturbation.

**Rule post-treatment.** By construction, there is some redundancy in the list of rules generated by the set of distinct paths  $\hat{\mathcal{P}}_{M,n,p_0}$ . The hyperrectangles associated with the paths extracted from a  $\Theta$ -random tree overlap, and so the corresponding rules are linearly dependent. Therefore a post-treatment to filter  $\hat{\mathcal{P}}_{M,n,p_0}$  is needed to remove redundancy and obtain a compact rule model. The general idea is straightforward: if the rule associated with the path  $\mathcal{P} \in \hat{\mathcal{P}}_{M,n,p_0}$  is a linear combination of rules associated with paths with a higher frequency in the forest, then  $\mathcal{P}$  is removed from  $\hat{\mathcal{P}}_{M,n,p_0}$ .

To illustrate the post-treatment, let the tree of Figure 1 be the  $\Theta_1$ -random tree grown in the forest. Since the paths of the first level of the tree,  $\mathcal{P}_1$  and  $\mathcal{P}_2$ , always occur in the same trees, we have  $\hat{p}_{M,n}(\mathcal{P}_1) = \hat{p}_{M,n}(\mathcal{P}_2)$ . If we assume these quantities to be greater than  $p_0$ , then  $\mathcal{P}_1$  and  $\mathcal{P}_2$  both belong to  $\hat{\mathcal{P}}_{M,n,p_0}$ . However, by construction,  $\mathcal{P}_1$  and  $\mathcal{P}_2$  are associated with the same rule, and we therefore enforce SIRUS to keep only  $\mathcal{P}_1$  in  $\hat{\mathcal{P}}_{M,n,p_0}$ . Each of the paths of the second level of the tree,  $\mathcal{P}_3$ ,  $\mathcal{P}_4$ ,  $\mathcal{P}_5$ , and  $\mathcal{P}_6$ , can occur in many different trees, and their associated  $\hat{p}_{M,n}$  are distinct (except in very specific cases). Assume for example that  $\hat{p}_{M,n}(\mathcal{P}_1) > \hat{p}_{M,n}(\mathcal{P}_4) > \hat{p}_{M,n}(\mathcal{P}_5) > \hat{p}_{M,n}(\mathcal{P}_3) > \hat{p}_{M,n}(\mathcal{P}_6) > p_0$ . Since  $\hat{g}_{n,\mathcal{P}_3}$  is a linear combination of  $\hat{g}_{n,\mathcal{P}_4}$  and  $\hat{g}_{n,\mathcal{P}_1}$ ,  $\mathcal{P}_3$  is removed. Similarly  $\mathcal{P}_6$  is redundant with  $\mathcal{P}_1$  and  $\mathcal{P}_5$ , and it is therefore removed. Finally, among the six paths of the tree, only  $\mathcal{P}_1$ ,  $\mathcal{P}_4$ , and  $\mathcal{P}_5$  are kept in the list  $\hat{\mathcal{P}}_{M,n,p_0}$ .

**Rule aggregation.** Now, the resulting small set of rules  $\hat{\mathcal{P}}_{M,n,p_0}$  is combined to form a simple, compact, and stable rule classification model. We simply average the set of elementary

rules  $\{\hat{g}_{n,\mathcal{P}} : \mathcal{P} \in \hat{\mathcal{P}}_{M,n,p_0}\}$  that have been selected in the first steps of SIRUS. The aggregated estimate  $\hat{\eta}_{M,n,p_0}(\mathbf{x})$  of  $\eta(\mathbf{x})$  is thus defined by

$$\hat{\eta}_{M,n,p_0}(\mathbf{x}) = \frac{1}{|\hat{\mathcal{P}}_{M,n,p_0}|} \sum_{\mathcal{P} \in \hat{\mathcal{P}}_{M,n,p_0}} \hat{g}_{n,\mathcal{P}}(\mathbf{x}). \quad (3.3)$$

Finally, the classification procedure assigns class 1 to an input  $\mathbf{x}$  if the aggregated estimate  $\hat{\eta}_{M,n,p_0}(\mathbf{x})$  is above a given threshold, and class 0 otherwise. In the introduction, we presented an example of a list of 7 rules for the Titanic dataset. In this case, for a new input  $\mathbf{x}$ ,  $\hat{\eta}_{M,n,p_0}(\mathbf{x})$  is simply the average of the output probability of survival  $p_s$  over the 7 selected rules.

In past works on rule ensemble models, such as RuleFit (Friedman and Popescu, 2008) and Node harvest (Meinshausen, 2010), rules are also extracted from a tree ensemble and then combined together through a regularized linear model. In our case, it happens that the parameter  $p_0$  alone is enough to control sparsity. Indeed, in our experiments, we observe that adding such linear model in the aggregation method hardly increases the accuracy and hardly reduces the size of the final rule set, while it can significantly reduce stability, add a set of coefficients that makes the model less straightforward to interpret, and requires more intensive computations. We refer to the experiments in Appendix A.3 for a comparison between  $\hat{\eta}_{M,n,p_0}$  defined as a simple average (3.3) versus a definition with a logistic regression.

**Categorical and numerical discrete variables.** For the sake of clarity, the description of SIRUS algorithm is limited to the case of numerical continuous variables. However, SIRUS can obviously handle numerical discrete and categorical data, as it is the case for random forests. On one hand, numerical discrete variables are left untouched since the number of possible split points is already finite, and the rule definition introduced for continuous variables also applies. On the other hand, we naturally extend the rule definition for categorical variables to “if  $X^{(1)}$  is *category 1 or 2* then *response* else *default response*”—see the Titanic dataset example in the introduction. Originally, categorical variables are efficiently handled in trees by transformation in ordered variables. Such ordering of categories is done with respect to the output mean for each category—see Breiman et al. (1984); Friedman et al. (2001), and we follow `ranger` implementation. Notice that trees are likely to often cut on categorical variables with a high number of categories, as highlighted in Strobl et al. (2006). Consequently, SIRUS is likely to output irrelevant rules associated to such categorical variables. Thus, it is best to discard categorical variables with a high number of categories, or transform them by regrouping categories or using one-hot-encoding before running SIRUS. Finally, note that ordinal variables (e.g.  $X^{(1)} \in \{\text{small, medium, big}\}$ ) are treated like categorical variables.

**Stability.** The three main properties to assess the interpretability of SIRUS are simplicity, stability, and predictivity, as already stated. On one hand, a measure of simplicity is naturally provided by the number of rules, and predictivity is given by the missclassification rate or the AUC. On the other hand, stability requires a more thorough discussion. In the statistical learning theory, stability refers to the stability of predictions (e.g., Vapnik, 1998). In particular, Rogers and Wagner (1978), Devroye and Wagner (1979), Bousquet and Elisseeff (2002), and Poggio et al. (2004) show that stability and predictive accuracy are closely connected. In our

case, we are more concerned by the stability of the internal structure of the model, and, to our knowledge, no general definition exists. So, we state the following tentative definition: a rule learning algorithm is stable if two independent estimations based on two independent samples result in two similar lists of rules. Thus, given a new sample  $\mathcal{D}'_n$  independent of  $\mathcal{D}_n$ , we define  $\hat{p}'_{M,n}(\mathcal{P})$  and the corresponding set of paths  $\hat{\mathcal{P}}'_{M,n,p_0}$  based on a modified random forest drawn with a parameter  $\Theta'$  independent of  $\Theta$ . Then, we measure the stability of SIRUS by the proportion of rules shared by the two sets  $\hat{\mathcal{P}}_{M,n,p_0}$  and  $\hat{\mathcal{P}}'_{M,n,p_0}$ , selected over these two runs of SIRUS on independent samples. We take advantage of a dissimilarity measure between two sets, the so-called Dice-Sorensen index, often used to assess the stability of variable selection methods (Chao et al., 2006; Zucknick et al., 2008; Boulesteix and Slawski, 2009; He and Yu, 2010; Alelyani et al., 2011). This index is defined by

$$\hat{S}_{M,n,p_0} = \frac{2|\hat{\mathcal{P}}_{M,n,p_0} \cap \hat{\mathcal{P}}'_{M,n,p_0}|}{|\hat{\mathcal{P}}_{M,n,p_0}| + |\hat{\mathcal{P}}'_{M,n,p_0}|} \quad (3.4)$$

with the convention  $0/0 = 1$ . This is a measure of stability taking values between 0 and 1: if the intersection between  $\hat{\mathcal{P}}_{M,n,p_0}$  and  $\hat{\mathcal{P}}'_{M,n,p_0}$  is empty, then  $\hat{S}_{M,n,p_0} = 0$ , while if  $\hat{\mathcal{P}}_{M,n,p_0} = \hat{\mathcal{P}}'_{M,n,p_0}$ , then  $\hat{S}_{M,n,p_0} = 1$ . Notice that it is possible to use other metrics to assess the distance between two finite sets (Zucknick et al., 2008): the Jaccard Index is another popular example. Although the stability values slightly vary with a different definition, both the asymptotic stability of SIRUS—see Section 4—and the empirical stability comparisons between algorithms—see Section 5—are insensitive to the stability metric choice.

## 4 Theoretical Analysis of Stability

Among the three minimum requirements for interpretability defined in Section 1, simplicity and predictivity are quite easily met for rule models (Cohen and Singer, 1999; Meinshausen, 2010; Letham et al., 2015). On the other hand, as Letham et al. (2015) recall, building a stable rule ensemble is challenging. Therefore the main goal of this section is to prove the asymptotic stability of SIRUS, i.e., provided that the sample size is large enough, SIRUS systematically outputs the same list of rules when run multiple times with independent samples. On the other hand, we also argue that existing tree-based rule algorithms are unstable by design.

In order to show the asymptotic stability of SIRUS, we first need to introduce formal definitions of the mathematical elements involved in the empirical algorithm. We additionally define the theoretical counterpart of SIRUS, an abstract procedure which is not based on the sample  $\mathcal{D}_n$ , but only on the unknown distribution  $\mathbf{P}_{\mathbf{X},Y}$ . Next, we will prove the stochastic convergence of SIRUS towards its theoretical counterpart. This means that the list of selected rules does not depend on the training data  $\mathcal{D}_n$ , but only on  $\mathbf{P}_{\mathbf{X},Y}$ , provided that the sample size is large enough. Therefore, the same list of rules is output when SIRUS is run multiple times on independent samples. This mathematical analysis highlights that the remarkable stable behavior of SIRUS in practice has theoretical groundings, and that the discretization of the cut values with the quantiles, as well as using random forests, are the cornerstones to stabilize rule models extracted from tree ensembles.

**Empirical algorithm.** First, we define the empirical CART-splitting criterion used to find the optimal split at each node of each tree of the forest. In our context of binary classification where the output  $Y \in \{0, 1\}$ , maximizing the so-called empirical CART-splitting criterion is equivalent to maximizing the criterion based on Gini impurity (see, e.g., [Biau and Scornet, 2016](#)). More precisely, at node  $H$  and for a cut performed along the  $j$ -th coordinate at the empirical  $r$ -th  $q$ -quantile  $\hat{q}_{n,r}^{(j)}$ , this criterion reads

$$\begin{aligned} L_n(H, \hat{q}_{n,r}^{(j)}) &\stackrel{\text{def}}{=} \frac{1}{N_n(H)} \sum_{i=1}^n (Y_i - \bar{Y}_H)^2 \mathbf{1}_{\mathbf{X}_i \in H} \\ &\quad - \frac{1}{N_n(H)} \sum_{i=1}^n (Y_i - \bar{Y}_{H_L} \mathbf{1}_{X_i^{(j)} < \hat{q}_{n,r}^{(j)}} - \bar{Y}_{H_R} \mathbf{1}_{X_i^{(j)} \geq \hat{q}_{n,r}^{(j)}})^2 \mathbf{1}_{\mathbf{X}_i \in H}, \end{aligned} \quad (4.1)$$

where  $\bar{Y}_H$  is the average of the  $Y_i$ 's such that  $\mathbf{X}_i \in H$ ,  $N_n(H)$  is the number of data points  $\mathbf{X}_i$  falling into  $H$ ,

$$H_L \stackrel{\text{def}}{=} \{\mathbf{x} \in H : \mathbf{x}^{(j)} < \hat{q}_{n,r}^{(j)}\}, \quad H_R \stackrel{\text{def}}{=} \{\mathbf{x} \in H : \mathbf{x}^{(j)} \geq \hat{q}_{n,r}^{(j)}\},$$

and for  $r \in \{1, \dots, q-1\}$  the empirical  $r$ -th  $q$ -quantile of  $\{X_1^{(j)}, \dots, X_n^{(j)}\}$  is defined by

$$\hat{q}_{n,r}^{(j)} = \inf \left\{ x \in \mathbb{R} : \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{X_i^{(j)} \leq x} \geq \frac{r}{q} \right\}. \quad (4.2)$$

Note that, for the ease of reading, (4.1) is defined for a tree built with the entire dataset  $\mathcal{D}_n$  without resampling. As it is often the case in the theoretical analysis of random forests, we assume throughout this section that the subsampling of  $a_n$  observations to build each tree is done without replacement to alleviate the mathematical analysis.

Recall that the rule selection is based on the probability  $p_n(\mathcal{P})$  that a  $\Theta$ -random tree of the forest contains a particular path  $\mathcal{P} \in \Pi$ , that is,

$$p_n(\mathcal{P}) = \mathbb{P}(\mathcal{P} \in T(\Theta, \mathcal{D}_n) | \mathcal{D}_n),$$

and that the Monte-Carlo estimate  $\hat{p}_{M,n}(\mathcal{P})$  of  $p_n(\mathcal{P})$  is directly computed using the random forest, and takes the form

$$\hat{p}_{M,n}(\mathcal{P}) = \frac{1}{M} \sum_{\ell=1}^M \mathbf{1}_{\mathcal{P} \in T(\Theta_\ell, \mathcal{D}_n)}.$$

Clearly,  $\hat{p}_{M,n}(\mathcal{P})$  is a good estimate of  $p_n(\mathcal{P})$  when  $M$  is large since, by the law of large numbers, conditional on  $\mathcal{D}_n$ ,

$$\lim_{M \rightarrow \infty} \hat{p}_{M,n}(\mathcal{P}) = p_n(\mathcal{P}) \quad \text{a.s.}$$

We also see that  $\hat{p}_{M,n}(\mathcal{P})$  is unbiased since  $\mathbb{E}[\hat{p}_{M,n}(\mathcal{P}) | \mathcal{D}_n] = p_n(\mathcal{P})$ .

**Theoretical algorithm.** Next, we define all theoretical counterparts of the empirical quantities involved in SIRUS, which do not depend on  $\mathcal{D}_n$  but only on the unknown distribution  $\mathbb{P}_{\mathbf{X},Y}$  of  $(\mathbf{X}, Y)$ . For a given integer  $q \geq 2$  and  $r \in \{1, \dots, q-1\}$ , the theoretical  $q$ -quantiles are defined by

$$q_r^{*(j)} = \inf\{x \in \mathbb{R} : \mathbb{P}(X^{(j)} \leq x) \geq \frac{r}{q}\},$$

i.e., the population version of  $\hat{q}_{n,r}^{(j)}$  defined in (4.2). Similarly, for a given hyperrectangle  $H \subseteq \mathbb{R}^p$ , we let the theoretical CART-splitting criterion be

$$\begin{aligned} L^*(H, q_r^{*(j)}) &= \mathbb{V}[Y|\mathbf{X} \in H] \\ &\quad - \mathbb{P}(X^{(j)} < q_r^{*(j)}|\mathbf{X} \in H) \times \mathbb{V}[Y|X^{(j)} < q_r^{*(j)}, \mathbf{X} \in H] \\ &\quad - \mathbb{P}(X^{(j)} \geq q_r^{*(j)}|\mathbf{X} \in H) \times \mathbb{V}[Y|X^{(j)} \geq q_r^{*(j)}, \mathbf{X} \in H]. \end{aligned}$$

Based on this criterion, we denote by  $T^*(\Theta)$  the list of all paths contained in the theoretical tree built with randomness  $\Theta$ , where splits are chosen to maximize the theoretical criterion  $L^*$  instead of the empirical one  $L_n$ , defined in (4.1). We stress again that the list  $T^*(\Theta)$  does not depend upon  $\mathcal{D}_n$  but only upon the unknown distribution of  $(\mathbf{X}, Y)$ . Next, we let  $p^*(\mathcal{P})$  be the theoretical counterpart of  $p_n(\mathcal{P})$ , that is

$$p^*(\mathcal{P}) = \mathbb{P}(\mathcal{P} \in T^*(\Theta)),$$

and finally define the theoretical set of selected paths  $\mathcal{P}_{p_0}^*$  by  $\{\mathcal{P} \in \Pi : p^*(\mathcal{P}) > p_0\}$  (with the same post-treatment as for the empirical procedure—see Section 3). Notice that, in the case where multiple splits have the same value of the theoretical CART-splitting criterion, one is randomly selected.

**Consistency of the path selection.** The construction of the rule ensemble model essentially relies on the path selection and on the estimates  $\hat{p}_{M,n}(\mathcal{P})$ ,  $\mathcal{P} \in \Pi$ . Therefore, our theoretical analysis first focuses on the asymptotic properties of those estimates in Theorem 1. Our consistency results hold under conditions on the subsampling rate  $a_n$  and the number of trees  $M_n$ , together with some assumptions on the distribution of the random vector  $\mathbf{X}$ . They are given below.

(A1) The subsampling rate  $a_n$  satisfies  $\lim_{n \rightarrow \infty} a_n = \infty$  and  $\lim_{n \rightarrow \infty} \frac{a_n}{n} = 0$ .

(A2) The number of trees  $M_n$  satisfies  $\lim_{n \rightarrow \infty} M_n = \infty$ .

(A3)  $\mathbf{X}$  has a strictly positive density  $f$  with respect to the Lebesgue measure. Furthermore, for all  $j \in \{1, \dots, p\}$ , the marginal density  $f^{(j)}$  of  $X^{(j)}$  is continuous, bounded, and strictly positive.

We can now state the consistency of the occurrence frequency of each possible path  $\mathcal{P} \in \Pi$  in the modified random forest.

**Theorem 1.** *If Assumptions (A1)-(A3) are satisfied, then, for all  $\mathcal{P} \in \Pi$ , we have*

$$\lim_{n \rightarrow \infty} \hat{p}_{M,n}(\mathcal{P}) = p^*(\mathcal{P}) \quad \text{in probability.}$$

**Stability.** The only source of randomness in the selection of the rules lies in the estimates  $\hat{p}_{M_n,n}(\mathcal{P})$ . Since Theorem 1 states the consistency of such an estimation, the path selection consistency follows, for all threshold values  $p_0$  that do not belong to the finite set  $\mathcal{U}^* = \{p^*(\mathcal{P}) : \mathcal{P} \in \Pi\}$  of all theoretical probabilities of appearance for each path  $\mathcal{P}$ . Indeed, if  $p_0 = p^*(\mathcal{P})$  for some  $\mathcal{P} \in \Pi$ , then  $\mathbb{P}(\hat{p}_{M_n,n}(\mathcal{P}) > p_0)$  does not necessarily converge to 0 and the path selection can be inconsistent. Then, we can deduce that SIRUS is asymptotically stable in the following Corollary 1.

**Corollary 1.** *Assume that Assumptions (A1)-(A3) are satisfied. Then, provided  $p_0 \in [0, 1] \setminus \mathcal{U}^*$ , we have*

$$\lim_{n \rightarrow \infty} \mathbb{P}(\hat{\mathcal{P}}_{M_n,n,p_0} = \mathcal{P}_{p_0}^*) = 1,$$

and then

$$\lim_{n \rightarrow \infty} \hat{S}_{M_n,n,p_0} = 1 \quad \text{in probability.}$$

**Competitors.** As we will discuss further in the experimental Section 5, CART, C5.0, Rule-Fit, and Node harvest are top competitors of SIRUS, which are also based on rule extraction from trees. However, these algorithms do not include a pre-processing step of discretization, which makes them unstable by design. To see this, we first adapt the definition of an extracted path without discretization as  $\mathcal{P} = \{(j_k, z_k, s_k), k = 1, \dots, d\}$ , where  $z_k \in \mathbb{R}$  is now the cutting value of the  $k$ -th split. For any rule algorithm, we also define  $\hat{S}_{M,n}$  as the proportion of rules shared between the output rule lists over two runs with two independent samples. Note that  $M = 1$  for CART and C5.0, and as already mentioned, it is possible to define a rule algorithm from CART, by extracting its nodes, as in C5.0. Thus, we obtain that for any tree-based rule algorithm,  $\hat{S}_{M,n} = 0$  almost surely. Indeed, since the input  $\mathbf{X}$  takes continuous values (Assumption (A3)) and decision trees can cut at the middle of two observations in all directions, the probability that a cutting value from the tree built with  $\mathcal{D}_n$  and one from the tree built with  $\mathcal{D}'_n$  are equal is null.

However, recall that in the experiments, we include a pre-processing discretization step to stabilize competitors and enable fair comparisons. With this modification, they reach a value of  $\hat{S}_{M,n} > 0$ , but still not in par with SIRUS. This shows that the high stability improvement of SIRUS does not only come from the discretization, but mainly from the rule selection procedure, based on the probability of the rule occurrence in a random tree.

**Proofs.** The proof of Theorem 1 is to be found in Appendix C. It is however interesting to give a sketch of the proof here. Corollary 1 is a direct consequence of Theorem 1, the full proof follows.

*Sketch of proof of Theorem 1.* The consistency is obtained by showing that  $\hat{p}_{M_n,n}(\mathcal{P})$  is asymptotically unbiased with a null variance. The result for the variance is quite straightforward since the variance of  $\hat{p}_{M_n,n}(\mathcal{P})$  can be broken into two terms: the variance generated by the Monte-Carlo randomization, which goes to 0 as the number of trees increases (Assumption (A2)), and the variance of  $p_n(\mathcal{P})$ . Following Mentch and Hooker (2016), since  $p_n(\mathcal{P})$  is a bagged estimate it can be seen as an infinite-order U-statistic, and a classic bound on the variance of U-statistics gives that  $\mathbb{V}[p_n(\mathcal{P})]$  converges to 0 if  $\lim_{n \rightarrow \infty} \frac{a_n}{n} = 0$ , which is true



by Assumption (A1). Next, proving that  $\hat{p}_{M_n,n}(\mathcal{P})$  is asymptotically unbiased requires to dive into the internal mechanisms of the random forest algorithm. To do this, we have to show that the CART-splitting criterion is consistent (Lemma 3) and asymptotically normal (Lemma 4) when cuts are limited to empirical quantiles (estimated on the same dataset) and the number of trees grows with  $n$ . When cuts are performed on the theoretical quantiles, the law of large numbers and the central limit theorem can be directly applied, so that the proof of Lemmas 3 and 4 boils down to showing that the difference between the empirical CART-splitting criterion evaluated at empirical and theoretical quantiles converges to 0 in probability fast enough. This is done in Lemma 2 thanks to Assumption (A3).  $\square$

*Proof of Corollary 1.* The first result is a consequence of Theorem 1 since

$$\mathbb{P}(\hat{\mathcal{P}}_{M_n,n,p_0} \neq \mathcal{P}_{p_0}^*) \leq \sum_{\mathcal{P} \in \Pi} \mathbb{P}(\hat{p}_{M_n,n}(\mathcal{P}) > p_0) \mathbb{1}_{p^*(\mathcal{P}) \leq p_0} + \mathbb{P}(\hat{p}_{M_n,n}(\mathcal{P}) \leq p_0) \mathbb{1}_{p^*(\mathcal{P}) > p_0}.$$

Next, we have

$$\hat{S}_{M_n,n,p_0} = \frac{2 \sum_{\mathcal{P} \in \Pi} \mathbb{1}_{\hat{p}_{M_n,n}(\mathcal{P}) > p_0} \mathbb{1}_{\hat{p}'_{M_n,n}(\mathcal{P}) > p_0}}{\sum_{\mathcal{P} \in \Pi} \mathbb{1}_{\hat{p}_{M_n,n}(\mathcal{P}) > p_0} + \mathbb{1}_{\hat{p}'_{M_n,n}(\mathcal{P}) > p_0}}.$$

Since  $p_0 \notin \mathcal{U}^*$ , we deduce from Theorem 1 and the continuous mapping theorem that, for all  $\mathcal{P} \in \Pi$ ,

$$\lim_{n \rightarrow \infty} \mathbb{1}_{\hat{p}_{M_n,n}(\mathcal{P}) > p_0} = \mathbb{1}_{p^*(\mathcal{P}) > p_0} \quad \text{in probability.}$$

Therefore,  $\lim_{n \rightarrow \infty} \hat{S}_{M_n,n,p_0} = 1$  in probability.  $\square$

## 5 Experiments

We begin this section by providing overall experimental settings. Next, we focus on a case study to illustrate SIRUS with an industrial process example: the semi-conductor manufacturing process SECOM data (Dua and Graff, 2017). In particular, it shows the excellent performance of SIRUS on real data in a noisy and high-dimensional setting. In Subsection 5.3, we use 19 UCI datasets (Dua and Graff, 2017) to perform extensive comparisons between SIRUS and its main competitors. We show that SIRUS produces much more stable rule lists, while preserving a predictive accuracy and computational complexity comparable to the top competitors. Finally, in Subsection 5.4, we detail the tuning procedure of the single hyperparameter  $p_0$ , along with a thorough discussion on the design of SIRUS. In particular, the cut limitations to the quantiles and the number of constraints in the selected rules are analyzed, and we also provide the stopping criterion for the number of trees.

### 5.1 Experiment Description

**Performance metrics.** We first introduce relevant metrics to assess the three interpretability properties in the experiments. By definition, the size (i.e., the simplicity) of the rule ensemble is the number of selected rules, i.e.,  $|\hat{\mathcal{P}}_{M,n,p_0}|$ . To measure the error, 1-AUC is used and

Dataset	Sample size	Total number of variables	Number of categorical variables
Authentication	1372	4	0
Breast Wisconsin	699	9	9
Credit Approval	690	15	9
Credit German	1000	20	13
Diabetes	768	8	0
Haberman	306	3	0
Heart C2	303	13	7
Heart H2	294	13	7
Heart Statlog	270	13	3
Hepatitis	155	19	0
Ionosphere	351	33	0
Kr vs Kp	3196	36	36
Liver Disorders	345	6	0
Mushrooms	8124	21	21
SECOM	1567	590	0
Sonar	208	60	0
Spambase	4601	57	0
Titanic	887	6	1
Vote	435	16	16
Wilt	4339	5	0

Table 1: Description of UCI datasets

estimated by 10-fold cross-validation (repeated 10 times for robustness and standard deviation estimates). With respect to stability, an independent dataset is not available for real data to compute  $\hat{S}_{M,n,p_0}$  as defined in (3.4) in the Section 3. Nonetheless, we can take advantage of the cross-validation process to compute a stability metric: the proportion of rules shared by two models built during the cross-validation, averaged over all possible pairs (Guidotti and Ruggieri, 2019).

**Datasets.** We have conducted experiments on the SECOM data, as well as 19 diverse public datasets from the UCI repository (Dua and Graff, 2017; data is described in Table 1). These experiments aim at illustrating the good behavior of SIRUS over its competitors in various settings. To compare stability of the different methods, data is discretized using the 10-empirical quantiles for each continuous variable and the same stability metric is used for all algorithm comparisons. For simplicity and predictivity metrics, we do not apply this pre-processing step of discretization, unless the algorithm only handles categorical data.

**Competitors.** For decision trees, we run both CART and C5.0, and trees are pruned to maximize their performance. Notice that, to enable simplicity and stability comparisons for CART, a list of rules is extracted from its nodes, as it is originally possible for C5.0. For rule algorithms based on greedy heuristics, we evaluate RIPPER, PART, and FOIL. Next,

for rule algorithms based on tree ensembles, we evaluate RuleFit and Node harvest. Note that categorical features are transformed in multiple binary variables as it is required by the two software implementations, and RuleFit is limited to rule predictors. For RuleFit, the lasso penalty is tuned by cross-validation as defined in [Friedman and Popescu \(2008\)](#). As advertised in [Meinshausen \(2010\)](#), Node harvest does not require parameter tuning by default, but it is also possible to add a regularization term to reduce the model size. We use the same tuning procedure as for SIRUS to maximize accuracy with the smallest possible model—see Subsection 5.4. Finally, for rule algorithms based on frequent pattern mining, we run the experiments for CBA and BRL. Note that we use default settings for BRL, since modifying its parameters does not significantly improve accuracy and can hurt stability. We use available R implementations: `rpart` ([Therneau and Atkinson, 2019](#), CART), `C50` ([Kuhn and Quinlan, 2020](#), C5.0), `RWeka` ([Hornik et al., 2009](#), RIPPER, PART), `arulesCBA` ([Johnson and Hahsler, 2020](#), FOIL, CBA), `pre` ([Fokkema, 2020](#), RuleFit), `nodeHarvest` ([Meinshausen, 2015](#), Node harvest), and `sbrl` ([Yang et al., 2017](#), BRL). We also use our R/C++ software implementation `sirus` ([Benard and Wright, 2020](#)) (available from CRAN), adapted from `ranger`, a fast random forest implementation ([Wright and Ziegler, 2017](#)). Besides, notice that for SIRUS experiments, we use the default settings of random forests well known for their excellent behavior, in particular `mtry` =  $\lfloor \frac{p}{3} \rfloor$ . We set  $q = 10$  quantiles and tune  $p_0$  as specified in Subsection 5.4.

## 5.2 Case Study: Manufacturing Process Data

SIRUS is run on a real manufacturing process of semi-conductors, the SECOM dataset ([Dua and Graff, 2017](#)). Data is collected from sensors and process measurement points to monitor the production line, resulting in 590 numeric variables. Each of the 1567 data points represents a single production entity associated with a pass or fail output (0/1) for in-house line testing. As it is often the case for a production process, the dataset is unbalanced and contains 104 fails, i.e., a failure rate  $p_f$  of 6.6%. We proceed to a simple pre-processing of the data: missing values (about 5% of the total) are replaced by the median.

Figure 2 shows predictivity versus the number of rules when  $p_0$  varies, with the optimal  $p_0$  displayed. Notice that the relation between  $p_0$  and the number of rules is monotone by construction, but also highly nonlinear. Therefore, we use the number of rules for the x-axis of Figure 2 to improve readability. The 1-AUC value is 0.30 for SIRUS (for the optimal  $p_0 = 0.04$ ), 0.29 for Breiman’s random forests, and 0.48 for a pruned CART tree. Thus, in that case, CART tree predicts no better than the random classifier, whereas SIRUS has a similar accuracy to random forests. The final model has 6 rules and a stability of 0.72, i.e., in average 4 to 5 rules are shared by 2 models built in a 10-fold cross-validation process, simulating data perturbation. By comparison, Node harvest outputs 36 rules with a value of 0.32 for 1-AUC.

Finally, the output of SIRUS may be displayed in the simple and interpretable form of Figure 3, the output in the R console of the package `sirus` for the SECOM data. Such a rule model enables to catch immediately how the most relevant variables impact failures. Among the 590 variables, 5 are enough to build a model as predictive as random forests, and such a selection is robust. Other rules alone may also be informative, but they do not add additional information to the model, since predictive accuracy is already minimal with the 6

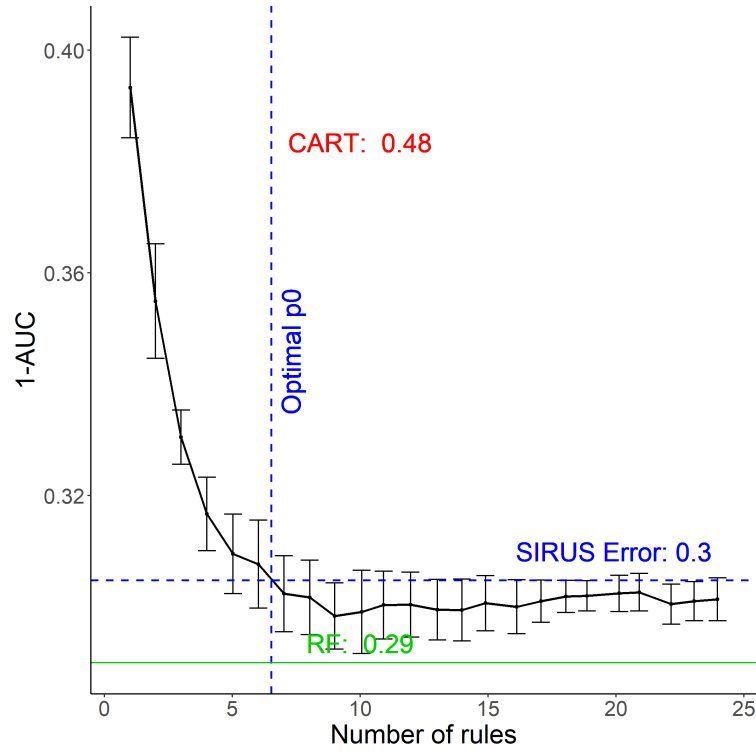


Figure 2: For the SECOM dataset, error (1-AUC) versus the number of rules when  $p_0$  varies, estimated via 10-fold cross-validation (averaged over 10 repetitions of the cross-validation). Errors for CART and random forests are reported for comparisons.

```
"Proportion of class 1 = 0.0664 - sample size n = 1567"
"if v60 < 5.51 then 0.0415 (n=1253) else 0.166 (n=314)"
"if v104 < -0.00868 then 0.0392 (n=1097) else 0.13 (n=470)"
"if v349 < 0.0356 then 0.0539 (n=1410) else 0.178 (n=157)"
"if v206 < 12.7 then 0.0539 (n=1410) else 0.178 (n=157)"
"if v65 < 26.1 then 0.0546 (n=1410) else 0.172 (n=157)"
"if v60 < 5.51 & v349 < 0.0356 then 0.0346 (n=1184) else 0.164 (n=383)"
```

Figure 3: List of rules output by our software `sirus` in the R console for the SECOM dataset.

selected rules. Then, production engineers should first focus on those 6 rules to investigate an improved setting of the production process. We insist that the stability of the output rule list is critical in practice. Indeed, the algorithm may be run multiple times during the analysis, eventually with an additional small new batch of data. The output rule list should be quite insensitive to such perturbation: domain experts are skeptical of unstable results, which are the symptoms of a partial and arbitrary modelling of the true phenomenon. SIRUS is stable, but it is not the case for decision trees or existing rule algorithms, as we show in the next subsection and illustrate in Appendix A.1.

### 5.3 Improvement over Competitors

Overall, we observe that SIRUS provides a high improvement of stability compared to state-of-the-art rule algorithms, while preserving the other properties. For the top competitors, experimental results are gathered in Table 2 for model size, Table 3 for stability, and Table 4 for predictive accuracy. Experiments for additional competitors are provided in Appendix A.2 in Tables 7, 8 and 9. Standard deviations are made negligible by averaging metrics over 10 repetitions of the cross-validation and are not displayed in the tables to increase readability.

Figure 4 provides an example for the dataset “Credit German” of the dependence between predictivity and the number of rules when  $p_0$  varies. In that case, the minimum of 1-AUC is about 0.25 for SIRUS, 0.20 for Breiman’s forests, and 0.29 for CART tree. For the chosen  $p_0$ , SIRUS returns a compact set of 22 rules and its stability is 0.66. Figure 5 provides another example of the good practical performance of SIRUS with the “Heart Statlog” dataset. Here, the predictivity of random forests is reached with 16 rules, with a stability of 0.83, i.e., about 13 rules are consistent between two different models built in a 10-fold cross-validation. Thus, the final models are simple, quite robust to data perturbation, and have a predictive accuracy close to random forests.

We can draw the following conclusions from the experimental comparisons with competitors, displayed in Tables 2, 3, and 4. SIRUS produces more stable and predictive rule lists than decision trees, for a comparable simplicity, but at the price of a higher computational complexity since many trees are grown. SIRUS produces much more stable and shorter rule lists than RuleFit and Node harvest, for a comparable accuracy and computational complexity. Classical rule algorithms exhibit similar properties as decision trees: a smaller computational complexity, but a high instability and a reduced predictivity. Finally, algorithms based on frequent pattern mining exhibit quite good stability properties, higher than for the other types of competitors. On the other hand, their predictive accuracy is worse than decision trees. Experiments in Tables 2, 3, and 4 show that SIRUS exhibits a high stability and predictivity improvement over these methods. Besides, simplicity varies across algorithms: CBA produces much longer rule lists than SIRUS, whereas BRL generates shorter models.

### 5.4 SIRUS Parameters

SIRUS relies on a single tuning hyperparameter: the selection threshold  $p_0$  involved in the definition of  $\hat{\mathcal{P}}_{M,n,p_0}$  to filter the most important rules, which therefore controls the simplicity of the model, and consequently also its accuracy and stability. On the other hand, SIRUS is

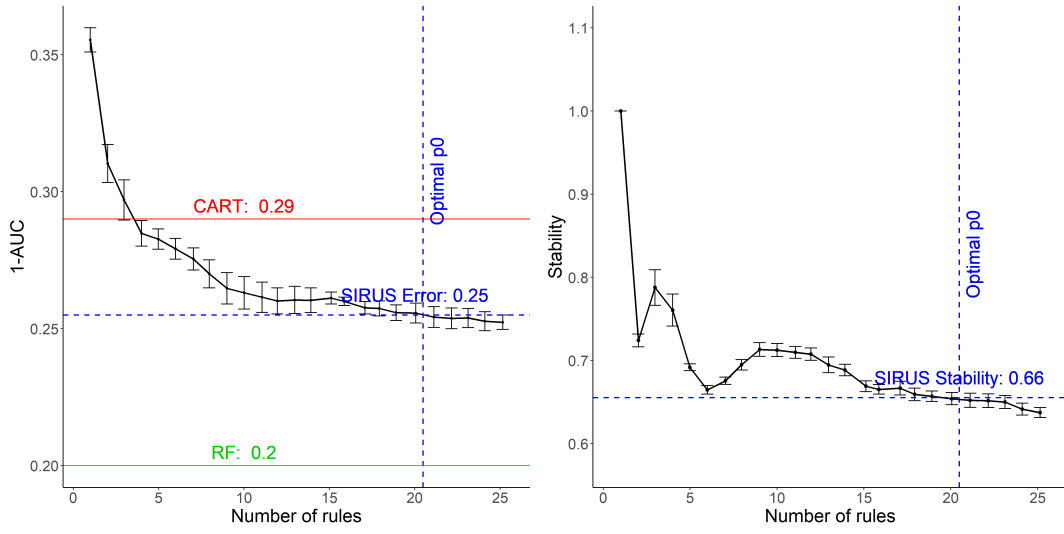


Figure 4: For the UCI dataset “Credit German”, 1-AUC (on the left) and stability (on the right) versus the number of rules when  $p_0$  varies, estimated via 10-fold cross-validation (results are averaged over 10 repetitions).

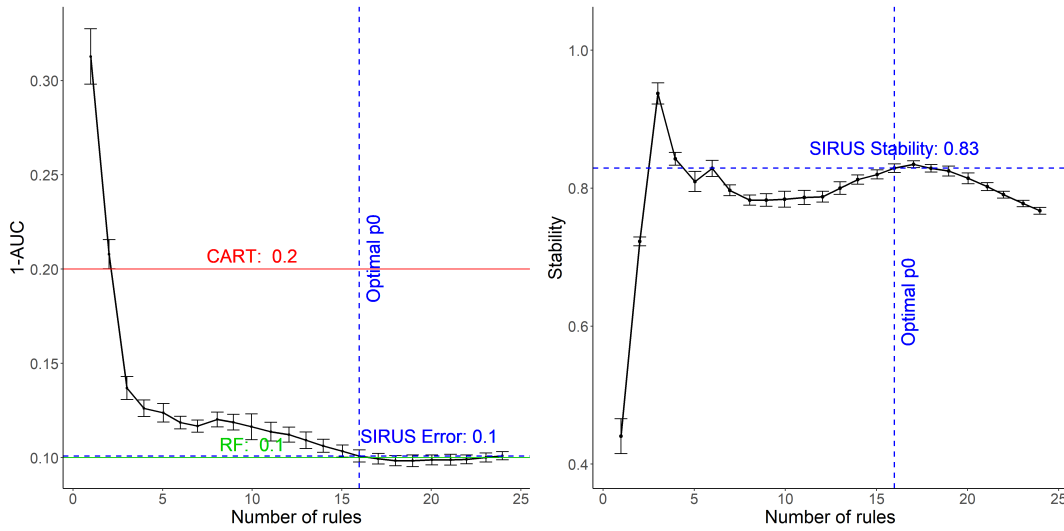


Figure 5: For the UCI dataset “Heart Statlog”, 1-AUC (on the left) and stability (on the right) versus the number of rules when  $p_0$  varies, estimated via 10-fold cross-validation (results are averaged over 10 repetitions).

	Decision tree	Classical rule learning	Frequent pattern mining		Tree ensemble		
Dataset	CART	RIPPER	CBA	BRL	RuleFit	Node harvest	SIRUS
Authentication	21	7	7	17	49	30	13
Breast Wisconsin	7	12	24	7	24	32	24
Credit Approval	5	4	55	4	15	27	16
Credit German	18	3	69	4	33	33	20
Diabetes	13	3	17	6	26	31	8
Haberman	2	1	2	2	3	17	5
Heart C2	10	3	34	4	23	36	20
Heart H2	5	2	29	3	12	24	12
Heart Statlog	10	3	27	4	22	35	16
Hepatitis	2	2	14	2	8	14	12
Ionosphere	4	4	38	4	20	35	15
Kr vs Kp	16	15	29	9	18	13	24
Liver Disorders	15	3	2	3	19	33	17
Mushrooms	4	8	25	11	10	22	23
Sonar	6	4	33	2	32	83	19
Spambase	14	16	126	16	68	60	21
Titanic	13	4	4	3	19	23	6
Vote	2	2	25	NA	12	10	7
Wilt	9	5	3	10	31	19	24

Table 2: Mean model size over a 10-fold cross-validation for UCI datasets. Results are averaged over 10 repetitions of the cross-validation.



	Decision tree	Classical rule learning	Frequent pattern mining		Tree ensemble		
Dataset	CART	RIPPER	CBA	BRL	RuleFit	Node harvest	SIRUS
Authentication	0.41	0.36	<b>0.87</b>	<b>0.86</b>	0.48	0.59	<b>0.81</b>
Breast Wisconsin	0.21	0.55	<b>0.80</b>	<b>0.78</b>	0.34	0.71	0.70
Credit Approval	0.52	0.32	0.43	0.52	0.25	0.23	<b>0.75</b>
Credit German	0.46	0.22	0.51	0.41	0.24	0.48	<b>0.66</b>
Diabetes	0.29	0.21	0.46	<b>0.73</b>	0.39	0.45	<b>0.81</b>
Haberman	<b>0.83</b>	0.09	<b>0.79</b>	0.50	0.46	0.52	0.65
Heart C2	0.25	0.35	0.38	0.60	0.39	0.49	<b>0.71</b>
Heart H2	0.46	0.27	0.52	<b>0.73</b>	0.29	0.29	<b>0.65</b>
Heart Statlog	0.30	0.41	0.41	<b>0.75</b>	0.35	0.48	<b>0.83</b>
Hepatitis	0.26	0.16	0.24	0.34	0.26	0.49	<b>0.68</b>
Ionosphere	<b>0.96</b>	0.39	0.13	0.70	0.17	0.33	0.69
Kr vs Kp	0.71	0.74	<b>0.84</b>	<b>0.80</b>	0.19	0.27	<b>0.87</b>
Liver Disorders	0.23	0.10	<b>0.91</b>	0.50	0.24	0.31	0.58
Mushrooms	<b>1</b>	0.84	<b>0.98</b>	0.80	0.69	0.48	0.86
Sonar	0.34	0.04	0.09	0.19	0.09	0.20	<b>0.55</b>
Spambase	0.49	0.10	0.46	<b>0.86</b>	0.28	0.66	<b>0.78</b>
Titanic	0.55	0.42	0.69	<b>0.88</b>	0.37	0.36	0.76
Vote	<b>1</b>	0.52	0.68	NA	0.21	0.30	0.75
Wilt	0.36	0.32	0.72	<b>0.94</b>	0.47	0.64	0.73
Average Rank	4.2	5.9	3.3	2.8	5.6	4.3	1.9
p-values	0.07	0.33	0.33	0.08	0.05	0.98	
<b>Final Rank</b>	<b>4</b>	<b>6</b>	<b>2</b>	<b>2</b>	<b>6</b>	<b>4</b>	<b>1</b>

Table 3: Mean stability over a 10-fold cross-validation for UCI datasets. Results are averaged over 10 repetitions of the cross-validation. Values within 10% of the maximum are displayed in bold. Algorithms are ranked with a Mann-Whitney-Wilcoxon test, the p-value with the previous performing algorithm determines the final rank (10%-level test).

	Black box	Decision tree	Classical rule learning	Frequent pattern mining		Tree ensemble		
Dataset	Random Forest	CART	RIPPER	CBA	BRL	RuleFit	Node harvest	SIRUS
Authentication	$10^{-4}$	0.02	0.02	0.14	0.009	<b><math>9.10^{-4}</math></b>	0.02	0.03
Breast Wisconsin	0.009	0.06	0.07	0.05	0.02	<b>0.01</b>	<b>0.01</b>	<b>0.01</b>
Credit Approval	0.07	0.14	0.15	0.14	0.11	<b>0.08</b>	<b>0.07</b>	0.09
Credit German	0.20	0.29	0.38	0.40	0.33	<b>0.23</b>	<b>0.26</b>	<b>0.25</b>
Diabetes	0.17	0.25	0.29	0.30	0.25	<b>0.18</b>	<b>0.19</b>	<b>0.19</b>
Haberman	0.31	0.48	0.39	0.50	0.43	<b>0.37</b>	<b>0.34</b>	<b>0.35</b>
Heart C2	0.10	0.19	0.23	0.17	0.23	0.12	0.12	<b>0.10</b>
Heart H2	0.11	0.23	0.24	0.24	0.16	<b>0.11</b>	<b>0.11</b>	<b>0.12</b>
Heart Statlog	0.10	0.20	0.21	0.17	0.22	0.12	0.12	<b>0.10</b>
Hepatitis	0.12	0.48	0.39	0.36	0.33	0.20	0.23	<b>0.17</b>
Ionosphere	0.02	0.11	0.12	0.13	0.10	<b>0.04</b>	0.07	0.07
Kr vs Kp	$9.10^{-4}$	0.02	<b>0.009</b>	0.05	0.01	<b>0.009</b>	0.04	0.04
Liver Disorders	0.23	0.33	0.35	0.48	0.44	<b>0.27</b>	0.30	0.35
Mushrooms	0	0.007	$3.10^{-5}$	$5.10^{-4}$	<b><math>2.10^{-5}</math></b>	$5.10^{-4}$	0.002	$6.10^{-4}$
Sonar	0.07	0.27	0.26	0.25	0.44	<b>0.12</b>	0.16	0.2
Spambase	0.01	0.11	0.08	0.12	0.05	<b>0.02</b>	0.04	0.07
Titanic	0.13	0.19	0.21	0.27	0.21	<b>0.14</b>	0.16	0.17
Vote	0.01	0.06	0.04	0.06	NA	<b>0.02</b>	<b>0.02</b>	<b>0.02</b>
Wilt	0.007	0.18	0.13	0.48	0.07	<b>0.02</b>	0.08	0.11
Average Rank		5	4.9	5.8	4.4	1.4	2.4	2.8
p-values		0.22	0.24	0.01	$6.10^{-3}$		0.01	0.34
<b>Final Rank</b>		<b>4</b>	<b>4</b>	<b>7</b>	<b>4</b>	<b>1</b>	<b>2</b>	<b>2</b>

Table 4: Model error (1-AUC) over a 10-fold cross-validation for UCI datasets. Results are averaged over 10 repetitions of the cross-validation. Values within 10% of the minimum are displayed in bold, random forest is put aside. Algorithms are ranked with a Mann-Whitney-Wilcoxon test, the p-value with the previous performing algorithm determines the final rank (10%-level test).

not very sensitive to the other parameters: the number of trees, the number of quantiles, and the tree depth. Therefore, they do not require fine tuning, and we simply set efficient default values as explained below.

**Tuning of SIRUS.** This parameter  $p_0$  should be set to optimize a tradeoff between the number of rules, stability, and accuracy. In practice, it is difficult to settle such a criterion, and we choose to optimize  $p_0$  to maximize the predictive accuracy with the smallest possible set of rules. To achieve this goal, we proceed as follows. The error 1-AUC is estimated by 10-fold cross-validation for a fine grid of  $p_0$  values, defined such that  $|\hat{\mathcal{P}}_{M,n,p_0}|$  varies from 1 to 25 rules. (We let 25 be an arbitrary upper bound on the maximum number of rules, considering that a bigger set is not readable anymore.) The randomization introduced by the partition of the dataset in the 10 folds of the cross-validation process has a significant impact on the variability of the size of the final model. Therefore, in order to get a robust estimation of  $p_0$ , the cross-validation is repeated multiple times (typically 10) and results are averaged. The standard deviation of the mean of 1-AUC is computed over these repetitions for each  $p_0$  of the grid search. We consider that all models within 2 standard deviations of the minimum of 1-AUC are not significantly less predictive than the optimal one. Thus, among these models, the one with the smallest number of rules is selected, i.e., the optimal  $p_0$  is shifted towards higher values to reduce the model size without decreasing predictivity—see Figures 4 and 5 for examples. This approach is very similar to the tuning procedure of the Lasso (Tibshirani, 1996).

**Number of trees.** The accuracy, stability, and computational cost of SIRUS increase with the number of trees  $M$ . Thus, we simply design a stopping criterion to grow the minimum number of trees which ensures that accuracy and stability are higher than 95% of their maximum asymptotic values with respect to  $M$  and conditionally on  $\mathcal{D}_n$ . We empirically observe that the stability requirement is met for a much higher number of trees than the accuracy requirement (about 10 times). Therefore, the stopping criterion is only based on stability. More precisely, we require that 95% of the rules are identical across two runs of SIRUS on a given dataset  $\mathcal{D}_n$  in average. Formally, the mean stability  $\mathbb{E}[\hat{S}_{M,n,p_0}|\mathcal{D}_n]$  measures the expected proportion of rules shared by two fits of SIRUS on  $\mathcal{D}_n$ , for fixed  $n$  (sample size),  $p_0$  (threshold), and  $M$  (number of trees). Thus, the stopping criterion takes the form  $1 - \mathbb{E}[\hat{S}_{M,n,p_0}|\mathcal{D}_n] < \alpha$ , with typically  $\alpha = 0.05$ .

There are two obstacles to operationalize this stopping criterion: its estimation and its dependence to  $p_0$ . We make two approximations to overcome these limitations and give empirical and theoretical evidence of their good practical behavior in Appendix B. First, Theorem 2 in Appendix B.2 provides an asymptotic equivalent with respect to  $M$  of  $1 - \mathbb{E}[\hat{S}_{M,n,p_0}|\mathcal{D}_n]$ , that we simply estimate by

$$\varepsilon_{M,n,p_0} = \frac{\sum_{\mathcal{P} \in \Pi} \Phi(Mp_0, M, \hat{p}_{M,n}(\mathcal{P}))(1 - \Phi(Mp_0, M, \hat{p}_{M,n}(\mathcal{P})))}{\sum_{\mathcal{P} \in \Pi} (1 - \Phi(Mp_0, M, \hat{p}_{M,n}(\mathcal{P})))},$$

where  $\Phi(Mp_0, M, p_n(\mathcal{P}))$  is the cdf of a binomial distribution with parameter  $p_n(\mathcal{P})$ ,  $M$  trials, evaluated at  $Mp_0$ . Secondly,  $\varepsilon_{M,n,p_0}$  depends on  $p_0$ , whose optimal value is unknown in the first step of SIRUS, when trees are grown. It turns out however that  $\varepsilon_{M,n,p_0}$  is not very sensitive

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

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

to  $p_0$ , as shown by the experiments in Appendix B.1. Consequently, our strategy is to simply average  $\varepsilon_{M,n,p_0}$  over a set  $\hat{V}_{M,n}$  of many possible values of  $p_0$  and use the resulting average as a gauge. These values are chosen to scan all possible path sets  $\hat{\mathcal{P}}_{M,n,p_0}$ , of size ranging from 1 to 50 paths. When a set of 50 paths is post-treated, its size reduces to around 25 paths (as explained in the previous paragraph, 25 is an arbitrarily threshold on the maximum number of rules above which a rule model is not readable anymore). In order to generate path sets of such sizes, values of  $p_0$  are chosen halfway between two distinct consecutive  $\hat{p}_{M,n}(\mathcal{P})$ ,  $\mathcal{P} \in \Pi$ , restricted to the highest 50 values. Thus, in the experiments, we utilize the following criterion to stop the growing of the forest, with typically  $\alpha = 0.05$ :

$$\operatorname{argmin}_M \left\{ \frac{1}{|\hat{V}_{M,n}|} \sum_{p_0 \in \hat{V}_{M,n}} \varepsilon_{M,n,p_0} < \alpha \right\}. \quad (5.1)$$

**Quantile discretization.** In the modified random forest grown in the first step of SIRUS, the split at each tree node is limited to the empirical  $q$ -quantiles of each component of  $\mathbf{X}$ , as described in Section 3. Thus, we check that this modification alone of the forest has little impact on its accuracy. Using the R package **ranger**, 1-AUC is estimated for each dataset with 10-fold cross-validation for  $q \in \{2, 5, 10, 20\}$ . We leave aside datasets with a majority of categorical variables, results are averaged over 10 repetitions of the cross-validation, and displayed in Table 5. Clearly, the decrease of accuracy generated by this discretization is small, and not very sensitive to  $q$ , provided that  $q$  is not too small. Thus,  $q = 10$  appears to be a good default choice from the experiments. In fact, the small impact of the discretization on the forest error is not surprising: with only  $p = 10$  input variables, the input space is split in a fine grid of  $10^{10}$  hyperrectangles for  $q = 10$  quantiles, providing a high flexibility to the modified random forest to identify local patterns.

**Tree depth.** When SIRUS is fit using fully grown trees, the final set of rules  $\hat{\mathcal{P}}_{M,n,p_0}$  contains almost exclusively rules made of one or two splits, and rarely of three splits. Although this

may appear surprising at first glance, this phenomenon is in fact expected. Indeed, rules made of multiple splits are extracted from deeper tree levels and are thus more sensitive to data perturbation by construction. This results in much smaller values of  $\hat{p}_{M,n}(\mathcal{P})$  for rules with a high number of splits, and then deletion from the final set of path through the threshold  $p_0$ :  $\hat{\mathcal{P}}_{M,n,p_0} = \{\mathcal{P} \in \Pi : \hat{p}_{M,n}(\mathcal{P}) > p_0\}$ . To illustrate this, let us consider the following typical example with  $p = 100$  input variables and  $q = 10$  quantiles. There are  $qp = 100 \times 10 = 10^3$  possible splits at the root node of a tree, and then  $2pq = 2 \cdot 10^3$  paths of one split. Since the left and right paths of one split at the root node are associated to the same rule, there are  $qp = 10^3$  distinct rules of one split, about  $(2qp)^2 \approx 10^6$  distinct rules of two splits, and about  $(2qp)^3 \approx 10^{10}$  distinct rules of three splits. Using only rules of one split is too restrictive since it generates a small model class (a thousand rules for 100 input variables) and does not handle variable interactions. On the other hand, rules of two splits are numerous (about one million) and thus provide a large flexibility to SIRUS. More importantly, since there are 10 billion rules of three splits, a stable selection of a few of them is clearly a difficult task, and such complex rules are naturally discarded by SIRUS.

In the software implementation `sirus`, the tree depth parameter `max.depth` is a modifiable input, set to 2 by default to reduce the computational cost while leaving the output list of rules almost untouched as explained above. We conduct experiments where SIRUS is run with a tree depth of 1, 2, and 3, and results are displayed in Table 6. Over the nineteen UCI datasets, rules of three splits appear in SIRUS rule list in only four cases, and a significant accuracy improvement over a tree depth of 2 occurs only once, for the ‘Mushrooms’ dataset. On the other hand, for all datasets except two, SIRUS outputs rules of two constraints, and predictivity is improved over a tree depth of 1 for half of the datasets. The Titanic example shows how the rule list is drastically simplified by limiting tree depth to 1, lowering the insights provided by SIRUS:

**Average survival rate  $p_s = 39\%$ .**

<b>if</b>	<b>sex is male</b>	<b>then</b>	$p_s = 19\%$	<b>else</b>	$p_s = 74\%$
<b>if</b>	$1^{st}$ or $2^{nd}$ class	<b>then</b>	$p_s = 56\%$	<b>else</b>	$p_s = 24\%$

This analysis of tree depth is not new. Indeed, both RuleFit (Friedman and Popescu, 2008) and Node harvest (Meinshausen, 2010) articles discuss the optimal tree depth for the rule extraction from a tree ensemble in their experiments. They both conclude that the optimal depth is 2. Hence, the same hard limit of 2 is used in Node harvest. RuleFit is slightly less restrictive: for each tree, its depth is randomly sampled with an exponential distribution concentrated on 2, but allowing few trees of depth 1, 3, and 4. We insist that they both reach such conclusion without considering stability issues, but only focusing on accuracy. Further considering stability properties consolidates that growing shallow trees is optimal for rule extraction from tree ensembles.

Dataset	SIRUS - depth = 1	SIRUS - depth = 2	SIRUS - depth = 3
Authentication	0.07	<b>0.03</b>	<b>0.03</b>
Breast Wisconsin	0.01	0.01	0.01
Credit Approval	0.11	<b>0.09</b>	<b>0.09</b>
Credit German	0.25	0.25	0.26
Diabetes	0.19	0.19	0.19
Haberman	0.35	0.35	0.35
Heart C2	0.11	<b>0.10</b>	0.11
Heart H2	0.12	0.12	0.12
Heart Statlog	0.11	<b>0.10</b>	<b>0.10</b>
Hepatitis	<b>0.15</b>	0.17	0.18
Ionosphere	0.07	0.07	0.07
Kr vs Kp	0.05	<b>0.04</b>	0.06
Liver Disorders	0.38	<b>0.35</b>	<b>0.35</b>
Mushrooms	$3.10^{-3}$	$6.10^{-4}$	<b><math>3.10^{-4}</math></b>
Sonar	0.19	0.2	0.2
Spambase	0.06	0.07	0.07
Titanic	0.19	<b>0.17</b>	<b>0.16</b>
Vote	0.02	0.02	0.02
Wilt	0.19	<b>0.11</b>	<b>0.11</b>

Table 6: SIRUS error (1-AUC) over a 10-fold cross-validation (averaged over 10 repetitions) when tree depth is limited to 1, 2 or 3. Values within 10% of the minimum are displayed in bold, except for datasets with no significant variations.

## 6 Conclusion

Interpretability of learning algorithms is required for applications involving critical decisions, for example the analysis of production processes in the manufacturing industry. Although interpretability does not have a precise definition, we argued that simplicity, stability, and predictivity are minimum requirements. In particular, decision trees and rule algorithms both combine a simple structure and a good accuracy for nonlinear data, and are thus considered as state-of-the-art interpretable algorithms. However, these methods are unstable with respect to data perturbation, which is a strong operational limitation. Therefore, we proposed a new rule algorithm for classification, SIRUS (Stable and Interpretable RULe Set), which takes the form of a short list of rules. We proved that SIRUS considerably improves stability over state-of-the-art algorithms, while preserving simplicity, accuracy, and computational complexity of top competitors. The principle of SIRUS is to extract rules from a random forest, based on their probability of occurrence in a random tree, and to stop the growing of the forest when the rule selection is converged. Thus, SIRUS inherits the computational complexity of random forests, and has only one tuning parameter. A software implementation, the R/C++ package `sirus` (Benard and Wright, 2020), is available from CRAN. Besides, we believe that the extension of SIRUS to regression is a promising future research direction: the main challenge is the construction of an appropriate rule aggregation framework to accurately estimate continuous outputs without hurting stability. Furthermore, although SIRUS has the ability to handle high-dimensional data, as illustrated with the SECOM dataset (590 inputs), specific variable selection strategies could be used to reduce the number of possible rules and then improve SIRUS performance.

## References

- R. Agrawal, T. Imieliński, and A. Swami. Mining association rules between sets of items in large databases. In *Proceedings of the 1993 ACM SIGMOD International Conference on Management of Data*, pages 207–216, New York, 1993. ACM.
- S. Alelyani, Z. Zhao, and H. Liu. A dilemma in assessing stability of feature selection algorithms. In *13th IEEE International Conference on High Performance Computing & Communication*, pages 701–707, Piscataway, 2011. IEEE.
- E. Angelino, N. Larus-Stone, D. Alabi, M. Seltzer, and C. Rudin. Learning certifiably optimal rule lists for categorical data. *The Journal of Machine Learning Research*, 18:8753–8830, 2017.
- C. Benard and M.N. Wright. *sirus: Stable and Interpretable RULe Set*, 2020. URL <https://CRAN.R-project.org/package=sirus>. R package version 0.2.1.
- G. Biau and E. Scornet. A random forest guided tour (with comments and a rejoinder by the author). *TEST*, 25:197–268, 2016.
- A.-L. Boulesteix and M. Slawski. Stability and aggregation of ranked gene lists. *Briefings in Bioinformatics*, 10:556–568, 2009.



- O. Bousquet and A. Elisseeff. Stability and generalization. *Journal of Machine Learning Research*, 2:499–526, 2002.
- L. Breiman. Bagging predictors. *Machine Learning*, 24:123–140, 1996.
- L. Breiman. Random forests. *Machine Learning*, 45:5–32, 2001a.
- L. Breiman. Statistical modeling: The two cultures (with comments and a rejoinder by the author). *Statistical Science*, 16:199–231, 2001b.
- L. Breiman. Setting up, using, and understanding random forests v3.1. Technical report, UC Berkeley, 2003a. URL [https://www.stat.berkeley.edu/~breiman/Using\\_random\\_forests\\_V3.1.pdf](https://www.stat.berkeley.edu/~breiman/Using_random_forests_V3.1.pdf).
- L. Breiman, J.H. Friedman, R.A. Olshen, and C.J. Stone. *Classification and Regression Trees*. Chapman & Hall/CRC, Boca Raton, 1984.
- A. Chao, R.L. Chazdon, R.K. Colwell, and T.-J. Shen. Abundance-based similarity indices and their estimation when there are unseen species in samples. *Biometrics*, 62:361–371, 2006.
- P. Clark and T. Niblett. The CN2 induction algorithm. *Machine Learning*, 3:261–283, 1989.
- W.W. Cohen. Fast effective rule induction. In *Proceedings of the Twelfth International Conference on Machine Learning*, pages 115–123, San Francisco, 1995. Morgan Kaufmann Publishers Inc.
- W.W. Cohen and Y. Singer. A simple, fast, and effective rule learner. In *Proceedings of the Sixteenth National Conference on Artificial Intelligence and Eleventh Conference on Innovative Applications of Artificial Intelligence*, pages 335–342, Palo Alto, 1999. AAAI Press.
- M. Cvitković, A.-S. Smith, and J. Pande. Asymptotic expansions of the hypergeometric function with two large parameters application to the partition function of a lattice gas in a field of traps. *Journal of Physics A: Mathematical and Theoretical*, 50:265206, 2017.
- K. Dembczyński, W. Kotłowski, and R. Słowiński. ENDER: A statistical framework for boosting decision rules. *Data Mining and Knowledge Discovery*, 21:52–90, 2010.
- L. Devroye and T. Wagner. Distribution-free inequalities for the deleted and holdout error estimates. *IEEE Transactions on Information Theory*, 25:202–207, 1979.
- F. Doshi-Velez and B. Kim. Towards a rigorous science of interpretable machine learning. *arXiv:1702.08608*, 2017.
- D. Dua and C. Graff. UCI machine learning repository, 2017. URL <http://archive.ics.uci.edu/ml>.
- M. Fokkema. Fitting prediction rule ensembles with R package pre. *Journal of Statistical Software*, 92:1–30, 2020.

- E. Frank and I.H. Witten. Generating accurate rule sets without global optimization. In *Proceedings of the Fifteenth International Conference on Machine Learning*, pages 144–151, San Francisco, 1998. Morgan Kaufmann Publishers Inc.
- A.A. Freitas. Comprehensible classification models: A position paper. *ACM SIGKDD Explorations Newsletter*, 15:1–10, 2014.
- J. Friedman, T. Hastie, and R. Tibshirani. *The Elements of Statistical Learning*, volume 1. Springer series in statistics New York, 2001.
- J.H. Friedman and B.E. Popescu. Importance sampled learning ensembles. Technical report, Stanford University, 2003.
- J.H. Friedman and B.E. Popescu. Predictive learning via rule ensembles. *The Annals of Applied Statistics*, 2:916–954, 2008.
- J. Fürnkranz and G. Widmer. Incremental reduced error pruning. In *Proceedings of the 11th International Conference on Machine Learning*, pages 70–77, San Francisco, 1994. Morgan Kaufmann Publishers Inc.
- R. Guidotti and S. Ruggieri. On the stability of interpretable models. In *International Joint Conference on Neural Networks*, pages 1–8, Piscataway, 2019. IEEE.
- R. Guidotti, A. Monreale, S. Ruggieri, F. Turini, F. Giannotti, and D. Pedreschi. A survey of methods for explaining black box models. *ACM Computing Surveys*, 51:1–42, 2018.
- Z. He and W. Yu. Stable feature selection for biomarker discovery. *Computational Biology and Chemistry*, 34:215–225, 2010.
- W. Hoeffding. A class of statistics with asymptotically normal distribution. *The Annals of Mathematical Statistics*, 19:293–325, 09 1948. doi: 10.1214/aoms/1177730196. URL <https://doi.org/10.1214/aoms/1177730196>.
- K. Hornik, C. Buchta, and A. Zeileis. Open-source machine learning: R meets Weka. *Computational Statistics*, 24:225–232, 2009.
- I. Johnson and M. Hahsler. *arulesCBA: Classification Based on Association Rules*, 2020. URL <https://CRAN.R-project.org/package=arulesCBA>. R package version 1.1.6.
- M. Kuhn and R. Quinlan. *C50: C5.0 Decision Trees and Rule-Based Models*, 2020. URL <https://CRAN.R-project.org/package=C50>. R package version 0.1.3.
- K. Kumbier, S. Basu, J.B. Brown, S. Celniker, and B. Yu. Refining interaction search through signed iterative random forests. *arXiv:1810.07287*, 2018.
- H. Lakkaraju, S.H. Bach, and J. Leskovec. Interpretable decision sets: A joint framework for description and prediction. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 1675–1684, New York, 2016. ACM.
- B. Letham. *Statistical learning for decision making: Interpretability, uncertainty, and inference*. PhD thesis, Massachusetts Institute of Technology, 2015.

- B. Letham, C. Rudin, T.H. McCormick, and D. Madigan. Interpretable classifiers using rules and bayesian analysis: Building a better stroke prediction model. *The Annals of Applied Statistics*, 9:1350–1371, 2015.
- Z.C. Lipton. The mythos of model interpretability. *arXiv:1606.03490*, 2016.
- B. Liu, W. Hsu, and Y. Ma. Integrating classification and association rule mining. In *Proceedings of the 14th International Conference on Knowledge Discovery and Data Mining*, volume 98, pages 80–86, New York, 1998. ACM.
- N. Meinshausen. Node harvest. *The Annals of Applied Statistics*, 4:2049–2072, 2010.
- N. Meinshausen. *Node harvest*, 2015. URL <https://CRAN.R-project.org/package=nodeHarvest>. R package version 0.7-3.
- L. Mentch and G. Hooker. Quantifying uncertainty in random forests via confidence intervals and hypothesis tests. *Journal of Machine Learning Research*, 17:841–881, 2016.
- R.S. Michalski. On the quasi-minimal solution of the general covering problem. In *Proceedings of the Fifth International Symposium on Information Processing*, pages 125–128, New York, 1969. ACM.
- W.J. Murdoch, C. Singh, K. Kumbier, R. Abbasi-Asl, and B. Yu. Interpretable machine learning: Definitions, methods, and applications. *arXiv:1901.04592*, 2019.
- T. Oates and D. Jensen. The effects of training set size on decision tree complexity. In *Proceedings of the 14th International Conference on Machine Learning*, pages 254–262, San Francisco, 1997. Morgan Kaufmann Publishers Inc.
- F.W.J. Olver, D.W. Lozier, R.F. Boisvert, and C.W. Clark. *NIST Handbook of Mathematical Functions Hardback and CD-ROM*. Cambridge University Press, 2010.
- C. Piech. Titanic dataset. <https://web.stanford.edu/class/archive/cs/cs109/cs109.1166/problem12.html>, 2016. Accessed: 2020-10-26.
- T. Poggio, R. Rifkin, S. Mukherjee, and P. Niyogi. General conditions for predictivity in learning theory. *Nature*, 428:419–422, 2004.
- J.R. Quinlan. Learning logical definitions from relations. *Machine learning*, 5:239–266, 1990.
- J.R. Quinlan. *C4.5: Programs for Machine Learning*. Morgan Kaufmann Publishers Inc., San Mateo, 1992.
- J.R. Quinlan and R.M. Cameron-Jones. Induction of logic programs: Foil and related systems. *New Generation Computing*, 13:287–312, 1995.
- M.T. Ribeiro, S. Singh, and C. Guestrin. Why should I trust you? Explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 1135–1144, New York, 2016. ACM.
- R.L. Rivest. Learning decision lists. *Machine Learning*, 2:229–246, 1987.

- W.H. Rogers and T.J. Wagner. A finite sample distribution-free performance bound for local discrimination rules. *The Annals of Statistics*, 6:506–514, 1978.
- C. Rudin. Please stop explaining black box models for high stakes decisions. *arXiv:1811.10154*, 2018.
- S. Rüping. *Learning interpretable models*. PhD thesis, Universität Dortmund, 2006.
- R.J. Serfling. *Approximation Theorems of Mathematical Statistics*, volume 162. John Wiley & Sons, 2009.
- C. Strobl, A.-L. Boulesteix, A. Zeileis, and T. Hothorn. Bias in random forest variable importance measures. In *Workshop on Statistical Modelling of Complex Systems*. Citeseer, 2006.
- T. Therneau and B. Atkinson. *rpart: Recursive Partitioning and Regression Trees*, 2019. URL <https://CRAN.R-project.org/package=rpart>. R package version 4.1-15.
- R. Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society. Series B (Methodological)*, pages 267–288, 1996.
- G. Tolomei, F. Silvestri, A. Haines, and M. Lalmas. Interpretable predictions of tree-based ensembles via actionable feature tweaking. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 465–474, New York, 2017. ACM.
- V. Vapnik. *Statistical Learning Theory*. Wiley, New York, 1998.
- S.M. Weiss and N. Indurkha. Lightweight rule induction. In *Proceedings of the Seventeenth International Conference on Machine Learning*, pages 1135–1142, San Francisco, 2000. Morgan Kaufmann Publishers Inc.
- M.N. Wright and A. Ziegler. ranger: A fast implementation of random forests for high dimensional data in C++ and R. *Journal of Statistical Software*, 77:1–17, 2017.
- H. Yang, C. Rudin, and M. Seltzer. Scalable Bayesian rule lists. In *Proceedings of the 34th International Conference on Machine Learning*, pages 3921–3930, Cambridge MA, 2017. JMLR.
- X. Yin and J. Han. CPAR: Classification based on predictive association rules. In *Proceedings of the 2003 SIAM International Conference on Data Mining*, pages 331–335, Philadelphia, 2003. SIAM.
- B. Yu. Stability. *Bernoulli*, 19:1484–1500, 2013.
- B. Yu and K. Kumbier. Three principles of data science: Predictability, computability, and stability (PCS). *arXiv:1901.08152*, 2019.
- M.J. Zaki, S. Parthasarathy, M. Ogihara, and W. Li. Parallel algorithms for discovery of association rules. *Data Mining and Knowledge Discovery*, 1:343–373, 1997.

M. Zucknick, S. Richardson, and E.A. Stronach. Comparing the characteristics of gene expression profiles derived by univariate and multivariate classification methods. *Statistical Applications in Genetics and Molecular Biology*, 7:1–34, 2008.

## A Additional Experiments

### A.1 Robustness Illustration

For the SECOM dataset used in the experimental Section 5 of the article, only three rule algorithms achieve the same predictivity as random forests: RuleFit, Node harvest, and SIRUS (1-AUC of 0.30, whereas CART and BRL are no better than the random classifier with an error of 1-AUC = 0.5). SIRUS produces a short and stable list of 6 rules, while RuleFit and Node harvest generate complex, long, and unstable rule lists. Rule algorithms based on tree ensembles are stochastic since they rely on the tree randomness  $\Theta_1, \dots, \Theta_M$ . Consequently, RuleFit and Node harvest output different rule lists when run multiple times on the same dataset. Such behavior is a strong limitation in practice, as domain experts become skeptical of the algorithm conclusions. On the other hand, SIRUS is built to have a robust rule extraction mechanism, and the same list of rules is output over multiple repetitions with the same data, as proved in Theorem 2 in the next Section.

To illustrate this, we run each algorithm twice on the SECOM dataset, and display the output models in Figure 6 for SIRUS, Figure 7 for Node harvest, and Figure 8 for RuleFit. We set the regularization parameter of Node harvest and SIRUS as explained in Subsection 5.3 of the article, to maximize accuracy with the smallest possible model: for Node harvest  $\lambda = 4$ , and for SIRUS  $p_0 = 0.04$ . RuleFit is tuned as defined in Friedman and Popescu (2008). Figures 7 and 8 show that the rule lists output by RuleFit and Node harvest are quite different across multiple runs with the exact same data, while SIRUS has the same output.

We also observe that for the same accuracy, RuleFit and Node harvest models are longer and more complex than SIRUS. In addition, rules are aggregated using weights to generate predictions. This is not the case for SIRUS, which simply averages the 6 output rules. Finally, we can also mention that manually increasing the regularization of Node harvest, to reduce the model size to 6 rules as in SIRUS, strongly hurts accuracy, which drops to 0.39.

### A.2 Additional Competitors

Additional experiments are provided to compare SIRUS to other competitors: C5.0 (Quinlan, 1992) (decision tree), PART (Frank and Witten, 1998), and FOIL (Quinlan and Cameron-Jones, 1995) (classical rule learning algorithms). Model size results are provided in Table 7, stability in Table 8, and error in Table 9. The stability and accuracy improvement of SIRUS is clear.

### A.3 Rule Aggregation

In Section 3 of the article,  $\hat{\eta}_{M,n,p_0}(\mathbf{x})$  (3.3) is a simple average of the set of rules, defined as

$$\hat{\eta}_{M,n,p_0}(\mathbf{x}) = \frac{1}{|\hat{\mathcal{P}}_{M,n,p_0}|} \sum_{\mathcal{P} \in \hat{\mathcal{P}}_{M,n,p_0}} \hat{g}_{n,\mathcal{P}}(\mathbf{x}). \quad (\text{A.1})$$

"Proportion of class 1 = 0.0664 - sample size n = 1567" "if v60 < 5.51 then 0.0415 (n=1253) else 0.166 (n=314)" "if v104 < -0.00868 then 0.0392 (n=1097) else 0.13 (n=470)" "if v349 < 0.0356 then 0.0539 (n=1410) else 0.178 (n=157)" "if v206 < 12.7 then 0.0539 (n=1410) else 0.178 (n=157)" "if v65 < 26.1 then 0.0546 (n=1410) else 0.172 (n=157)" "if v60 < 5.51 & v349 < 0.0356 then 0.0346 (n=1184) else 0.164 (n=383)"
"Proportion of class 1 = 0.0664 - sample size n = 1567" "if v60 < 5.51 then 0.0415 (n=1253) else 0.166 (n=314)" "if v104 < -0.00868 then 0.0392 (n=1097) else 0.13 (n=470)" "if v349 < 0.0356 then 0.0539 (n=1410) else 0.178 (n=157)" "if v206 < 12.7 then 0.0539 (n=1410) else 0.178 (n=157)" "if v65 < 26.1 then 0.0546 (n=1410) else 0.172 (n=157)" "if v60 < 5.51 & v349 < 0.0356 then 0.0346 (n=1184) else 0.164 (n=383)"

Figure 6: The two lists of rules output by two runs of SIRUS for the SECOM dataset.

Dataset	C5.0	PART	FOIL	SIRUS
Authentication	11	8	20	13
Breast Wisconsin	5	10	41	24
Credit Approval	9	32	40	16
Credit German	22	68	101	22
Diabetes	12	7	36	8
Haberman	2	2	4	5
Heart C2	10	20	31	20
Heart H2	4	15	29	12
Heart Statlog	10	18	28	15
Hepatitis	7	8	14	12
Ionosphere	9	6	28	15
Kr vs Kp	11	21	24	24
Liver Disorders	14	7	2	17
Mushrooms	7	9	14	23
Sonar	10	6	20	19
Spambase	29	46	73	21
Titanic	7	15	17	6
Vote	5	7	19	7
Wilt	10	8	10	24

Table 7: Mean model size over a 10-fold cross-validation for UCI datasets (averaged over 10 repetitions).



```

"if v122 > 16 & v52 > 189 then 0.6 (n=10, weight=0.38)"
"if v122 > 16 & v481 < 56 then 0.545 (n=11, weight=0.12)"
"if v511 > 105 & v206 > 14.1 then 0.692 (n=13, weight=0.074)"
"if v65 > 30.7 & v116 < 722 then 0.643 (n=14, weight=0.207)"
"if v60 > 8.36 & v442 > 1.11 then 0.571 (n=14, weight=0.036)"
"if v122 > 16 & v481 > 56 then 0.143 (n=14, weight=0.12)"
"if v122 > 16 & v52 < 189 then 0.133 (n=15, weight=0.38)"
"if v104 < -0.00865 & v435 > 19.7 then 0.263 (n=19, weight=0.027)"
"if v60 < 4.96 & v122 > 16 then 0.304 (n=23, weight=0.027)"
"if v65 > 30.7 & v449 < 0.207 then 0.522 (n=23, weight=0.071)"
"if v60 > 8.41 & v521 < 1.3 then 0.462 (n=26, weight=0.08)"
"if v60 < 4.97 & v572 < 1.21 then 0.258 (n=31, weight=0.019)"
"if v60 < 8.04 & v65 > 33.3 then 0.294 (n=34, weight=0.223)"
"if v60 < 8.14 & v349 > 0.0443 then 0.257 (n=35, weight=0.129)"
"if v60 > 4.96 & v342 > 4.13 then 0.436 (n=39, weight=0.027)"
"if v60 > 8.14 & v334 > 6.76 then 0.475 (n=40, weight=0.352)"
"if v65 > 30.7 & v449 > 0.207 then 0.14 (n=43, weight=0.071)"
"if v65 > 30.7 & v116 > 722 then 0.173 (n=52, weight=0.207)"
"if v60 > 4.95 & v334 > 6.76 then 0.389 (n=54, weight=0.019)"
"if v104 > -0.00865 & v542 > 11.4 then 0.305 (n=82, weight=0.027)"
"if v511 > 105 & v206 < 14.1 then 0.108 (n=83, weight=0.074)"
"if v60 > 8.41 & v588 > 0.0161 then 0.292 (n=106, weight=0.106)"
"if v60 > 8.41 & v588 < 0.0161 then 0.106 (n=132, weight=0.106)"
"if v60 > 8.14 & v334 < 6.76 then 0.132 (n=204, weight=0.129)"
"if v60 > 8.04 & v334 < 6.76 then 0.136 (n=206, weight=0.223)"
"if v60 > 8.41 & v521 > 1.3 then 0.156 (n=212, weight=0.08)"
"if v60 > 8.41 & v171 < 0.971 then 0.165 (n=224, weight=0.036)"
"if v511 < 105 & v60 > 5.49 then 0.156 (n=269, weight=0.074)"
"if v60 > 4.97 & v334 < 6.76 then 0.125 (n=288, weight=0.019)"
"if v60 > 4.96 & v342 < 4.13 then 0.132 (n=304, weight=0.027)"
"if v104 > -0.00865 & v542 < 11.4 then 0.093 (n=388, weight=0.027)"
"if v104 < -0.00865 & v435 < 19.7 then 0.035 (n=1078, weight=0.027)"
"if v60 < 4.97 & v572 > 1.21 then 0.033 (n=1194, weight=0.019)"
"if v60 < 4.96 & v122 < 16 then 0.033 (n=1201, weight=0.027)"
"if v511 < 105 & v60 < 5.49 then 0.037 (n=1202, weight=0.074)"
"if v60 < 8.04 & v65 < 33.3 then 0.037 (n=1287, weight=0.223)"
"if v60 < 8.14 & v349 < 0.0443 then 0.038 (n=1288, weight=0.129)"
"if v60 < 8.41 & v122 < 16 then 0.039 (n=1304, weight=0.222)"
"if v65 < 30.7 & v122 < 16 then 0.053 (n=1476, weight=0.278)"

```

```

"if v104 > -0.00665 & v334 > 7.37 then 0.667 (n=12, weight=0.019)"
"if v65 > 30.7 & v457 < 8.81 then 0.692 (n=13, weight=0.322)"
"if v407 > 14 then 0.385 (n=13, weight=0.017)"
"if v60 < 8.08 & v430 > 10.4 then 0.333 (n=15, weight=0.171)"
"if v60 > 8.08 & v170 > 0.584 then 0.562 (n=16, weight=0.124)"
"if v17 > 10.8 then 0.278 (n=18, weight=0.019)"
"if v60 > 8.08 & v521 < 1.09 then 0.526 (n=19, weight=0.012)"
"if v60 < 8.08 & v122 > 16 then 0.292 (n=24, weight=0.096)"
"if v66 < 36.8 & v122 > 16 then 0.32 (n=25, weight=0.066)"
"if v133 < 2.21 then 0.296 (n=27, weight=0.129)"
"if v60 < 5.02 & v572 < 1.2 then 0.258 (n=31, weight=0.078)"
"if v66 > 36.4 & v342 > 3.19 then 0.455 (n=33, weight=0.066)"
"if v60 < 9.02 & v349 > 0.0438 then 0.25 (n=40, weight=0.124)"
"if v60 > 8.08 & v334 > 6.76 then 0.475 (n=40, weight=0.378)"
"if v60 < 8.94 & v65 > 31.7 then 0.255 (n=47, weight=0.124)"
"if v65 > 30.7 & v457 > 8.81 then 0.17 (n=53, weight=0.322)"
"if v66 > 36.4 & v342 < 3.19 then 0.104 (n=77, weight=0.066)"
"if v60 > 5.02 & v478 > 8.1 then 0.314 (n=86, weight=0.077)"
"if v65 < 30.7 & v60 > 10.9 then 0.189 (n=201, weight=0.176)"
"if v60 > 8.14 & v334 < 6.76 then 0.132 (n=204, weight=0.124)"
"if v60 > 8.08 & v334 < 6.76 then 0.137 (n=205, weight=0.145)"
"if v60 > 8.08 & v521 > 1.09 then 0.164 (n=226, weight=0.012)"
"if v60 > 6.54 & v334 < 6.76 then 0.131 (n=229, weight=0.109)"
"if v104 > -0.00665 & v334 < 7.37 then 0.13 (n=230, weight=0.019)"
"if v60 > 8.04 & v170 < 0.584 then 0.165 (n=230, weight=0.124)"
"if v60 > 5.02 & v478 < 8.1 then 0.115 (n=253, weight=0.077)"
"if v60 < 5.02 & v572 > 1.2 then 0.033 (n=1197, weight=0.078)"
"if v60 < 8.08 & v65 < 31.6 then 0.035 (n=1275, weight=0.124)"
"if v60 < 8.04 & v349 < 0.0438 then 0.037 (n=1282, weight=0.124)"
"if v60 < 6.54 & v430 < 10.4 then 0.039 (n=1283, weight=0.109)"
"if v60 < 8.08 & v122 < 16 then 0.039 (n=1298, weight=0.096)"
"if v65 < 30.7 & v60 < 10.9 then 0.037 (n=1300, weight=0.176)"
"if v104 < -0.00665 & v17 < 10.8 then 0.047 (n=1307, weight=0.019)"
"if v60 < 8.08 & v430 < 10.4 then 0.04 (n=1307, weight=0.062)"
"if v66 < 36.4 & v122 < 16 then 0.051 (n=1432, weight=0.066)"
"if v133 > 2.21 & v65 < 30.7 then 0.053 (n=1474, weight=0.129)"
"if v65 < 30.7 & v407 < 14 then 0.054 (n=1488, weight=0.017)"

```

Figure 7: The two lists of rules output by two runs of Node harvest for the SECOM dataset.

rule	coefficient	description
(Intercept)	-1.304499863	1
rule616	-0.400252692	v60 <= 4.97 & v105 > -0.0019 & v424 <= 108.6217
rule26	-0.399674943	v349 <= 0.0385 & v60 <= 8.3918 & v64 <= 17.6454
rule496	-0.265685341	v60 <= 0.8045 & v101 <= 5e-04 & v568 <= 0.0896
rule441	-0.260900593	v60 <= 7.8264 & v583 <= 0.5011 & v350 <= 0.049
rule314	-0.258822916	v22 <= -5512.5 & v472 <= 30.7812
rule508	-0.190299769	v511 <= 95.5975 & v101 <= 5e-04 & v153 <= 0.7523
rule43	-0.177421075	v60 <= 8.3918 & v349 <= 0.0342 & v139 <= 90.8
rule97	-0.134937737	v511 <= 95.3413 & v153 <= 0.7523 & v196 <= 0.361
rule444	-0.117968967	v104 <= -0.0087 & v34 <= 9.1637
rule368	-0.087452989	v104 <= -0.0079 & v153 <= 0.8257
rule395	-0.084409096	v65 <= 25.1618 & v60 <= 9.5927 & v438 <= 7.9865
rule628	-0.084144279	v130 <= 0.0946 & v350 <= 0.0611 & v361 <= 0.0036
rule86	-0.023078885	v125 <= 16.05 & v60 <= 4.9555 & v303 <= 0.45
rule362	-0.003972723	v104 <= -0.0087 & v436 <= 10.2733 & v350 <= 0.0595

rule	coefficient	description
(Intercept)	0.178336422	1
rule97	-0.523012600	v349 <= 0.0421 & v511 <= 200.823 & v60 <= 8.1445
rule282	-0.463529803	v511 <= 65.1163 & v153 <= 0.8257 & v197 <= 14.43
rule606	-0.338103339	v432 <= 99.2163 & v438 <= 7.1906 & v65 <= 30.5136
rule496	-0.297717157	v250 <= 0.0034 & v65 <= 25.1618 & v125 <= 16.05
rule289	-0.278210742	v456 <= 3.7084 & v288 <= 0.3448 & v555 <= 0.852
rule674	-0.272413104	v153 <= 0.7377 & v125 <= 16.04
rule404	-0.266285107	v60 <= 4.9382 & v303 <= 0.4304 & v105 > -0.0017
rule556	-0.261565996	v250 <= 8e-04 & v130 <= 0.0946 & v361 <= 0.0029
rule600	-0.258720261	v512 <= 708.5714 & v558 <= 2.9289 & v65 <= 30.68
rule500	-0.245999282	v22 <= -5394.25 & v438 <= 7.3595
rule461	-0.197524877	v22 <= -5581
rule197	-0.166101239	v104 <= -0.0087 & v301 <= 0.121 & v34 <= 9.7836
rule635	-0.157494908	v334 <= 6.6293 & v366 <= 0.013
rule92	-0.156029423	v349 <= 0.0362 & v511 <= 95.5975 & v438 <= 5.1928
rule130	-0.145965819	v104 <= -0.0087 & v299 <= 0.1024 & v41 > 14
rule140	-0.121309793	v349 <= 0.0369 & v472 <= 21.8646 & v60 <= 4.9991
rule84	-0.120009890	v60 <= 5.4718 & v104 <= -0.0067 & v526 <= 7.5026
rule171	-0.085220151	v334 <= 5.4943
rule595	-0.079847068	v34 <= 8.5891
rule571	-0.078349545	v60 <= 1.6018 & v526 <= 8.8106
rule36	-0.067557526	v60 <= 8.3918 & v511 <= 80.4829 & v349 <= 0.0441
rule361	-0.053981777	v349 <= 0.0369 & v511 <= 167.2026 & v334 <= 6.1301
rule368	-0.041471470	v65 <= 31.4709 & v60 <= 9.8518 & v168 <= 1.1
rule636	-0.037163161	v334 <= 6.6293 & v366 <= 0.013 & v34 <= 9.088
rule150	-0.032344454	v60 <= 4.92 & v349 <= 0.0437 & v288 <= 0.3456
rule448	-0.014851459	v130 <= 0.1892 & v350 <= 0.0595
rule521	-0.014601179	v60 <= 8.1445 & v511 <= 204.5307 & v65 <= 31.2182
rule177	0.013482768	v334 > 5.4943 & v104 > -0.0068
rule335	-0.012690307	v317 > 5.9229 & v520 <= 26.109 & v350 <= 0.0495
rule542	-0.005889676	v349 <= 0.0326 & v115 <= 7e-04 & v438 <= 7.1856

Figure 8: The two lists of rules output by two runs of RuleFit for the SECOM dataset.

Dataset	C5.0	PART	FOIL	SIRUS
Authentication	0.44	0.43	<b>0.81</b>	<b>0.81</b>
Breast Wisconsin	0.17	0.49	0.36	<b>0.70</b>
Credit Approval	0.18	0.31	0.17	<b>0.75</b>
Credit German	0.03	0.16	0.11	<b>0.65</b>
Diabetes	0.07	0.15	0.18	<b>0.81</b>
Haberman	0.28	0.25	<b>0.64</b>	<b>0.65</b>
Heart C2	0.09	0.15	0.16	<b>0.71</b>
Heart H2	0.32	0.31	0.39	<b>0.65</b>
Heart Statlog	0.11	0.15	0.15	<b>0.82</b>
Hepatitis	0.10	0.15	0.05	<b>0.68</b>
Ionosphere	0.24	0.13	0.07	<b>0.69</b>
Kr vs Kp	0.65	0.51	0.85	<b>0.87</b>
Liver Disorders	0.05	0.07	<b>0.69</b>	0.58
Mushrooms	0.79	0.78	<b>0.93</b>	<b>0.86</b>
Sonar	0.06	0.06	0.04	<b>0.55</b>
Spambase	0.08	0.08	0.11	<b>0.78</b>
Titanic	0.49	0.27	<b>0.77</b>	<b>0.76</b>
Vote	<b>0.67</b>	0.40	0.39	<b>0.75</b>
Wilt	0.34	0.37	0.48	<b>0.73</b>

Table 8: Mean stability over a 10-fold cross-validation for UCI datasets (averaged over 10 repetitions). Values within 10% of the maximum are displayed in bold.

Dataset	C5.0	PART	FOIL	SIRUS
Authentication	0.02	<b>0.01</b>	0.08	0.03
Breast Wisconsin	0.06	0.07	0.08	<b>0.01</b>
Credit Approval	0.15	0.17	0.15	<b>0.09</b>
Credit German	0.37	0.36	0.41	<b>0.25</b>
Diabetes	0.28	0.30	0.28	<b>0.19</b>
Haberman	0.46	0.42	0.50	<b>0.35</b>
Heart C2	0.20	0.23	0.19	<b>0.10</b>
Heart H2	0.23	0.23	0.23	<b>0.12</b>
Heart Statlog	0.21	0.24	0.20	<b>0.10</b>
Hepatitis	0.34	0.34	0.39	<b>0.17</b>
Ionosphere	0.10	0.10	0.13	<b>0.07</b>
Kr vs Kp	<b>0.006</b>	0.008	0.02	0.04
Liver Disorders	<b>0.34</b>	<b>0.38</b>	0.50	<b>0.35</b>
Mushrooms	0.001	0	<b><math>6.10^{-5}</math></b>	$6.10^{-4}$
Sonar	0.26	0.26	0.26	<b>0.2</b>
Spambase	<b>0.07</b>	<b>0.07</b>	0.12	<b>0.07</b>
Titanic	0.20	0.20	0.25	<b>0.17</b>
Vote	0.04	0.05	0.05	<b>0.02</b>
Wilt	0.15	0.17	0.46	<b>0.11</b>

Table 9: Model error (1-AUC) over a 10-fold cross-validation for UCI datasets (averaged over 10 repetitions). Values within 10% of the minimum are displayed in bold.

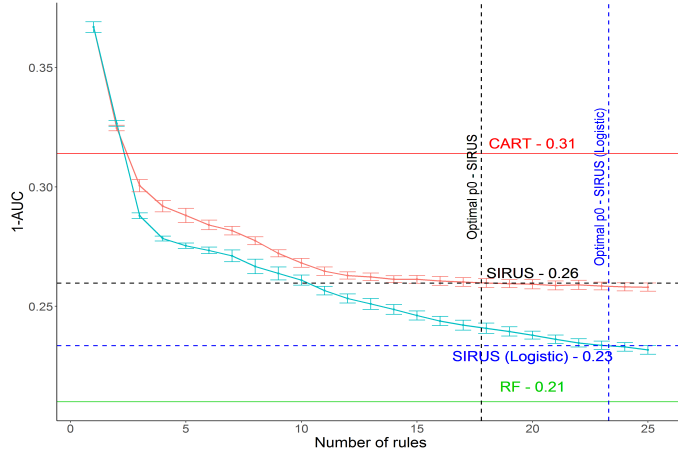


Figure 9: For the UCI dataset “Credit German”, 1-AUC versus the number of rules when  $p_0$  varies, estimated via 10-fold cross-validation (repeated 30 times) for two different methods of rule aggregation: the rule average (A.1) in red and a logistic regression (A.2) in blue.

To tackle our binary classification problem, a natural approach would be to use a logistic regression and define

$$\ln \left( \frac{\hat{\eta}_{M,n,p_0}(\mathbf{x})}{1 - \hat{\eta}_{M,n,p_0}(\mathbf{x})} \right) = \sum_{\mathcal{P} \in \mathcal{P}_{M,n,p_0}} \beta_{\mathcal{P}} \hat{g}_{n,\mathcal{P}}(\mathbf{x}), \quad (\text{A.2})$$

where the coefficients  $\beta_{\mathcal{P}}$  have to be estimated. To illustrate the performance of the logistic regression (A.2), we consider again the UCI dataset, “Credit German”. We augment the previous results from Figure 4 (in Section 5 of the article) with the logistic regression error in Figure 9. One can observe that the predictive accuracy is slightly improved but it comes at the price of an additional set of coefficients that can be hard to interpret (some can be negative), and an increased computational cost. Notice that categorical variables are one-hot-encoded in this example.

## B Stopping Criterion for the Number of Trees $M$

We recall that the definition of the stopping criterion (5.1) of the forest growing is provided in Section 5 of the main article. First, we provide three groups of experiments to show its good empirical efficiency. In the second subsection, we provide theoretical properties of the stopping criterion.

### B.1 Experiments

The following experiments on the UCI datasets show the good empirical performance of the stopping criterion (5.1). Recall that the goal of this criterion is to determine the minimum number of trees  $M$  ensuring that two independent fits of SIRUS on the same dataset result

Dataset	Mean stability
Haberman	0.950 (0.01)
Diabetes	0.950 (0.007)
Heart Statlog	0.954 (0.007)
Liver Disorders	0.951 (0.006)
Heart C2	0.955 (0.009)
Heart H2	0.952 (0.009)
Credit German	0.950 (0.008)
Credit Approval	0.941 (0.02)
Ionosphere	0.950 (0.009)

Table 10: Values of  $\hat{S}_{M,n,p_0}$  averaged over  $p_0 \in \hat{V}_{M,n}$  when the stopping criterion (5.1) is used to set  $M$ , for UCI datasets. Results are averaged over 10 repetitions and standard deviations are displayed in parentheses.

in two lists of rules with an overlap of 95% in average. This is checked with a first batch of experiments—see next paragraph. Secondly, the stopping criterion (5.1) does not consider the optimal  $p_0$ , unknown when trees are grown in the first step of SIRUS. Then, another batch of experiments is run to show that the stability approximation  $1 - \varepsilon_{M,n,p_0}$  is quite insensitive to  $p_0$ . Finally, a last batch of experiments provides examples of the number of trees grown when SIRUS is fit. Notice that for these experiments, categorical variables are one-hot-encoded.

**Experiments 1.** For each dataset, the following procedure is applied. SIRUS is run a first time using criterion (5.1) to stop the number of trees. This initial run provides the optimal number of trees  $M$  as well as the set  $\hat{V}_{M,n}$  of possible  $p_0$ . Then, SIRUS is fit twice independently using the precomputed number of trees  $M$ . For each  $p_0 \in \hat{V}_{M,n}$ , the stability metric  $\hat{S}_{M,n,p_0}$  (with  $\mathcal{D}'_n = \mathcal{D}_n$ ) is computed over the two resulting lists of rules. Finally  $\hat{S}_{M,n,p_0}$  is averaged across all  $p_0$  values in  $\hat{V}_{M,n}$ . This procedure is repeated 10 times: results are averaged and presented in Table 10, with standard deviations in parentheses. Across the considered datasets, resulting values range from 0.941 to 0.955, and are thus close to 0.95 as expected by construction of criterion (5.1).

**Experiments 2.** The second type of experiments illustrates that  $\varepsilon_{M,n,p_0}$  is quite insensitive to  $p_0$  when  $M$  is set with criterion (5.1). For the “Credit German” dataset, we fit SIRUS and then compute  $1 - \varepsilon_{M,n,p_0}$  for each  $p_0 \in \hat{V}_{M,n}$ . Results are displayed in Figure 10.  $1 - \varepsilon_{M,n,p_0}$  ranges from 0.90 to 1, where the extreme values are reached for  $p_0$  corresponding to very small number of rules, which are not of interest when  $p_0$  is selected to maximize predictive accuracy. Thus,  $1 - \varepsilon_{M,n,p_0}$  is quite concentrated around 0.95 when  $p_0$  varies.

**Experiments 3.** Finally, we display in Table 11 the optimal number of trees when the growing of SIRUS is stopped using criterion (5.1). It ranges from 4220 to 20 650 trees. In Breiman’s forests, the number of trees above which the accuracy cannot be significantly improved is typically 10 times lower. However SIRUS grows shallow trees, and is thus not computationally

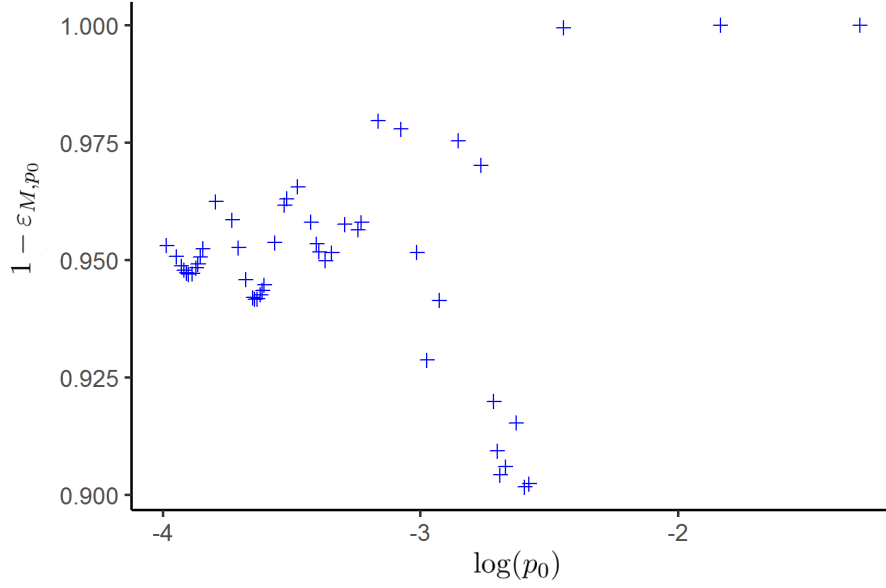


Figure 10: For the UCI dataset “Credit German”,  $1 - \varepsilon_{M,n,p_0}$  for a sequence of  $p_0 \in \hat{V}_{M,p_0}$  corresponding to final models ranging from 1 to about 25 rules.

more demanding than random forests overall.

## B.2 Theoretical Properties

We emphasize that growing more trees does not improve predictive accuracy or stability with respect to data perturbation for a fixed sample size  $n$ . Indeed, the instability of the rule selection is generated by the variance of the estimates  $\hat{p}_{M,n}(\mathcal{P})$ ,  $\mathcal{P} \in \Pi$ . Upon noting that we have two sources of randomness— $\Theta$  and  $\mathcal{D}_n$ —, the law of total variance shows that  $\mathbb{V}[\hat{p}_{M,n}(\mathcal{P})]$  can be broken down into two terms: the variance generated by the Monte Carlo randomness  $\Theta$  on the one hand, and the sampling variance on the other hand. In fact, equation (C.3) in the proof of Theorem 1 below reveals that

$$\mathbb{V}[\hat{p}_{M,n}(\mathcal{P})] = \frac{1}{M} \mathbb{E}[p_n(\mathcal{P})](1 - \mathbb{E}[p_n(\mathcal{P})]) + (1 - \frac{1}{M}) \mathbb{V}[p_n(\mathcal{P})].$$

The stopping criterion (5.1) ensures that the first term becomes negligible as  $M \rightarrow \infty$ , so that  $\mathbb{V}[\hat{p}_{M,n}(\mathcal{P})]$  reduces to the sampling variance  $\mathbb{V}[p_n(\mathcal{P})]$ , which is independent of  $M$ . Therefore, stability with respect to data perturbation cannot be further improved by increasing the number of trees. Additionally, the trees are only involved in the selection of the paths. For a given set of paths  $\hat{\mathcal{P}}_{M,n,p_0}$ , the construction of the final aggregated estimate  $\hat{\eta}_{M,n,p_0}$  (see Section 3 of the article) is independent of the forest. Thus, if further increasing the number of trees does not impact the path selection, neither it improves the predictive accuracy.

Next, Theorem 2 states that conditionally on  $\mathcal{D}_n$  and with  $\mathcal{D}'_n = \mathcal{D}_n$ ,  $\hat{S}_{M,n,p_0}$  should be close to 1, and also provides an asymptotic approximation of  $\mathbb{E}[\hat{S}_{M,n,p_0} | \mathcal{D}_n]$  for large values of

Dataset	Nb of trees (sd)
Haberman	10 920 (877)
Diabetes	18 830 (1538)
Heart Statlog	7840 (994)
Liver Disorders	14 650 (1242)
Heart C2	6840 (1270)
Heart H2	4220 (529)
Credit German	7940 (672)
Credit Approval	20 650 (8460)
Ionosphere	7320 (487)

Table 11: Number of trees  $M$  determined by the stopping criterion (5.1) for UCI datasets. Results are averaged over 10 repetitions and standard deviations are displayed in parentheses.

the number of trees  $M$ , which quantifies the influence of  $M$  on the mean stability, conditional on  $\mathcal{D}_n$ . We let  $\mathcal{U}_n \stackrel{\text{def}}{=} \{p_n(\mathcal{P}) : \mathcal{P} \in \Pi\}$  be the empirical counterpart of  $\mathcal{U}^*$ .

**Theorem 2.** *If  $p_0 \in [0, 1] \setminus \mathcal{U}_n$  and  $\mathcal{D}'_n = \mathcal{D}_n$ , then, conditional on  $\mathcal{D}_n$ , we have*

$$\lim_{M \rightarrow \infty} \hat{S}_{M,n,p_0} = 1 \quad \text{in probability.} \quad (\text{B.1})$$

*In addition, for all  $p_0 < \max \mathcal{U}_n$ ,*

$$1 - \mathbb{E}[\hat{S}_{M,n,p_0} | \mathcal{D}_n] \underset{M \rightarrow \infty}{\sim} \sum_{\mathcal{P} \in \Pi} \frac{\Phi(Mp_0, M, p_n(\mathcal{P}))(1 - \Phi(Mp_0, M, p_n(\mathcal{P})))}{\frac{1}{2} \sum_{\mathcal{P}' \in \Pi} \mathbb{1}_{p_n(\mathcal{P}') > p_0} + \mathbb{1}_{p_n(\mathcal{P}') > p_0 - \rho_n(\mathcal{P}, \mathcal{P}') \frac{\sigma_n(\mathcal{P}')}{\sigma_n(\mathcal{P})} (p_0 - p_n(\mathcal{P}))}},$$

where  $\Phi(Mp_0, M, p_n(\mathcal{P}))$  is the cdf of a binomial distribution with parameter  $p_n(\mathcal{P})$ ,  $M$  trials, evaluated at  $Mp_0$ , and, for all  $\mathcal{P}, \mathcal{P}' \in \Pi$ ,

$$\sigma_n(\mathcal{P}) = \sqrt{p_n(\mathcal{P})(1 - p_n(\mathcal{P}))},$$

and

$$\rho_n(\mathcal{P}, \mathcal{P}') = \frac{\text{Cov}(\mathbb{1}_{\mathcal{P} \in T(\Theta, \mathcal{D}_n)}, \mathbb{1}_{\mathcal{P}' \in T(\Theta, \mathcal{D}_n)} | \mathcal{D}_n)}{\sigma_n(\mathcal{P})\sigma_n(\mathcal{P}')}.$$

The proof of Theorem 2 is to be found in Section D. The equivalent provided in Theorem 2 is defined when the sets of rules  $\hat{\mathcal{P}}_{M,n,p_0}$  and  $\hat{\mathcal{P}}'_{M,n,p_0}$  are not post-treated. It considerably simplifies the analysis of the asymptotic behavior of  $\mathbb{E}[\hat{S}_{M,n,p_0} | \mathcal{D}_n]$ . Since the post-treatment is deterministic, this operation is not an additional source of instability. Then, if the estimation of the rule set without post-treatment is stable, it is also the case when the post-treatment is added. Finally, despite its apparent complexity, the asymptotic approximation of  $1 - \mathbb{E}[\hat{S}_{M,n,p_0} | \mathcal{D}_n]$  can be easily estimated, and an efficient stopping criterion for the number of trees is therefore derived in (5.1).



## C Proof of Theorem 1

We recall Assumptions (A1)-(A3) and Theorem 1 for the sake of clarity.

(A1) The subsampling rate  $a_n$  satisfies  $\lim_{n \rightarrow \infty} a_n = \infty$  and  $\lim_{n \rightarrow \infty} \frac{a_n}{n} = 0$ .

(A2) The number of trees  $M_n$  satisfies  $\lim_{n \rightarrow \infty} M_n = \infty$ .

(A3)  $\mathbf{X}$  has a strictly positive density  $f$  with respect to the Lebesgue measure. Furthermore, for all  $j \in \{1, \dots, p\}$ , the marginal density  $f^{(j)}$  of  $X^{(j)}$  is continuous, bounded, and strictly positive.

**Theorem 1.** *If Assumptions (A1)-(A3) are satisfied, then, for all  $\mathcal{P} \in \Pi$ , we have*

$$\lim_{n \rightarrow \infty} \hat{p}_{M_n, n}(\mathcal{P}) = p^*(\mathcal{P}) \quad \text{in probability.}$$

First, we prove Theorem 1 for a path of one split. The proof is extended for a path of two splits in the next subsection and follows the same steps. Finally, the proof can be easily extended to a path of any depth  $d \in \mathbb{N}^*$  by recursion.

### C.1 Proof of Theorem 1 for a path of one split

We consider  $\mathcal{P}_1 = \{(j_1, r_1, s_1)\}$  a path of one split, where  $j_1 \in \{1, \dots, p\}$ ,  $r_1 \in \{1, \dots, q-1\}$ , and  $s_1 \in \{L, R\}$ . We assume throughout that Assumptions (A1)-(A3) are satisfied.

Before proving Theorem 1, we state five lemmas (Lemma 1 to Lemma 5). Their proof can be found in the Subsection C.3. Lemma 1 is a preliminary technical result used to state both Lemmas 2 and 4 - case (b).

**Lemma 1.** *Let  $\mathbf{X}$  be a random variable distributed on  $\mathbb{R}^p$  such that Assumptions (A1) and (A3) are satisfied. Then, for all  $j \in \{1, \dots, p\}$  and all  $r \in \{1, \dots, q-1\}$ , we have*

$$\lim_{n \rightarrow \infty} \sqrt{a_n} \mathbb{P}(q_r^{*(j)} \leq X^{(j)} < \hat{q}_{n,r}^{(j)}) = 0$$

and

$$\lim_{n \rightarrow \infty} \sqrt{a_n} \mathbb{P}(\hat{q}_{n,r}^{(j)} \leq X^{(j)} < q_r^{*(j)}) = 0.$$

Lemma 2 is used to prove both consistency (Lemma 3) and convergence rate (Lemma 4) of the CART-splitting criterion when the root node of the tree is cut at an empirical quantile. Lemma 5 is an intermediate result to prove Theorem 1.

**Lemma 2.** *If Assumptions (A1) and (A3) are satisfied, then for all  $j \in \{1, \dots, p\}$ , all  $r \in \{1, \dots, q-1\}$ , and all  $H \subseteq \mathbb{R}^p$  such that  $\mathbb{P}(\mathbf{X} \in H, X^{(j)} < q_r^{*(j)}) > 0$  and  $\mathbb{P}(\mathbf{X} \in H, X^{(j)} \geq q_r^{*(j)}) > 0$ , we have*

$$\lim_{n \rightarrow \infty} \sqrt{a_n} (L_{a_n}(H, \hat{q}_{n,r}^{(j)}) - L_{a_n}(H, q_r^{*(j)})) = 0 \quad \text{in probability.}$$



**Lemma 3.** *If Assumptions (A1) and (A3) are satisfied, then for all  $j \in \{1, \dots, p\}$ , all  $r \in \{1, \dots, q-1\}$ , and all  $H \subseteq \mathbb{R}^p$  such that  $\mathbb{P}(\mathbf{X} \in H, X^{(j)} < q_r^{*(j)}) > 0$  and  $\mathbb{P}(\mathbf{X} \in H, X^{(j)} \geq q_r^{*(j)}) > 0$ , we have*

$$\lim_{n \rightarrow \infty} L_{a_n}(H, \hat{q}_{n,r}^{(j)}) = L^*(H, q_r^{*(j)}) \quad \text{in probability.}$$

When splitting a node, if the theoretical CART-splitting criterion has multiple maxima, one is randomly selected. This random selection follows a discrete probability law, which is not necessarily uniform and is based on  $\mathbb{P}_{\mathbf{X}, Y}$  as specified in Definition 1. In order to derive the limit of the probability that a given split occurs in a  $\Theta$ -random tree in the empirical algorithm, one needs to assess the convergence rate of the empirical CART-splitting criterion when it has multiple maxima.

**Lemma 4.** *Consider that Assumptions (A1) and (A3) are satisfied. Let  $\mathcal{C}_1 \subset \{1, \dots, p\} \times \{1, \dots, q-1\}$  be a set of splits of cardinality  $c_1 \geq 2$ , such that, for all  $(j, r) \in \mathcal{C}_1$ ,  $L^*(\mathbb{R}^p, q_r^{*(j)}) \stackrel{\text{def}}{=} L_{\mathcal{C}_1}^*$ , i.e., the theoretical CART-splitting criterion is constant for all splits in  $\mathcal{C}_1$ . Let  $(j_1, r_1) \in \mathcal{C}_1$  and let  $\mathbf{L}_{n, \mathcal{P}_1}^{(\mathcal{C}_1)}$  be a random vector where each component is the difference between the empirical CART-splitting criterion for the splits  $(j, r) \in \mathcal{C}_1 \setminus (j_1, r_1)$  and  $(j_1, r_1)$ , that is*

$$\mathbf{L}_{n, \mathcal{P}_1}^{(\mathcal{C}_1)} = \left( L_{a_n}(\mathbb{R}^p, \hat{q}_{n,r}^{(j)}) - L_{a_n}(\mathbb{R}^p, \hat{q}_{n,r_1}^{(j_1)}) \right)_{(j,r) \in \mathcal{C}_1 \setminus (j_1, r_1)}.$$

(a) *If  $L_{\mathcal{C}_1}^* > 0$ , then we have*

$$\sqrt{a_n} \mathbf{L}_{n, \mathcal{P}_1}^{(\mathcal{C}_1)} \xrightarrow[n \rightarrow \infty]{\mathcal{D}} \mathcal{N}(0, \Sigma),$$

where, for all  $(j, r), (j', r') \in \mathcal{C}_1 \setminus (j_1, r_1)$ , each element of the covariance matrix  $\Sigma$  is defined by  $\Sigma_{(j,r),(j',r')} = \text{Cov}[Z_{j,r}, Z_{j',r'}]$ , with

$$\begin{aligned} Z_{j,r} = & \left( Y - \mathbb{E}[Y|X^{(j_1)} < q_{r_1}^{*(j_1)}] \mathbb{1}_{X^{(j_1)} < q_{r_1}^{*(j_1)}} \right. \\ & \left. - \mathbb{E}[Y|X^{(j_1)} \geq q_{r_1}^{*(j_1)}] \mathbb{1}_{X^{(j_1)} \geq q_{r_1}^{*(j_1)}} \right)^2 \\ & - \left( Y - \mathbb{E}[Y|X^{(j)} < q_r^{*(j)}] \mathbb{1}_{X^{(j)} < q_r^{*(j)}} - \mathbb{E}[Y|X^{(j)} \geq q_r^{*(j)}] \mathbb{1}_{X^{(j)} \geq q_r^{*(j)}} \right)^2. \end{aligned}$$

Besides, for all  $(j, r) \in \mathcal{C}_1$ ,  $\mathbb{V}[Z_{j,r}] > 0$ .

(b) *If  $L_{\mathcal{C}_1}^* = 0$ , then we have*

$$a_n \mathbf{L}_{n, \mathcal{P}_1}^{(\mathcal{C}_1)} \xrightarrow[n \rightarrow \infty]{\mathcal{D}} h_{\mathcal{P}_1}(\mathbf{V}),$$

where  $\mathbf{V}$  is a Gaussian vector of covariance matrix  $\text{Cov}[\mathbf{Z}]$ . If  $\mathcal{C}_1$  is explicitly written  $\mathcal{C}_1 = \{(j_k, r_k)\}_{k=1, \dots, c_1}$ ,  $\mathbf{Z}$  is defined, for  $k \in \{1, \dots, c_1\}$ , by

$$\begin{aligned} Z_{2k-1} &= \frac{1}{\sqrt{p_{L,k}}} (Y - \mathbb{E}[Y]) \mathbb{1}_{X^{(j_k)} < q_{r_k}^{*(j_k)}} \\ Z_{2k} &= \frac{1}{\sqrt{p_{R,k}}} (Y - \mathbb{E}[Y]) \mathbb{1}_{X^{(j_k)} \geq q_{r_k}^{*(j_k)}}, \end{aligned}$$

where  $p_{L,k} = \mathbb{P}(X^{(j_k)} < q_{r_k}^{*(j_k)})$ ,  $p_{R,k} = \mathbb{P}(X^{(j_k)} \geq q_{r_k}^{*(j_k)})$ , and  $h_{\mathcal{P}_1}$  is a multivariate quadratic form defined as

$$h_{\mathcal{P}_1} : \begin{pmatrix} x_1 \\ \vdots \\ x_{2c_1} \end{pmatrix} \rightarrow \begin{pmatrix} x_3^2 + x_4^2 - x_1^2 - x_2^2 \\ \vdots \\ x_{2k-1}^2 + x_{2k}^2 - x_1^2 - x_2^2 \\ \vdots \\ x_{2c_1-1}^2 + x_{2c_1}^2 - x_1^2 - x_2^2 \end{pmatrix}.$$

Besides, the variance of each component of  $h_{\mathcal{P}_1}(\mathbf{V})$  is strictly positive.

**Definition 1** (Theoretical splitting procedure). Let  $\theta_1^{(V)}$  be the set of eligible variables to split the root node of a theoretical random tree. The set of best theoretical cuts at the root node is defined as

$$\mathcal{C}_1^*(\theta_1^{(V)}) = \underset{(j,r) \in \theta_1^{(V)} \times \{1, \dots, q-1\}}{\operatorname{argmax}} L^*(\mathbb{R}^p, q_r^{*(j)}).$$

If  $\mathcal{C}_1^*(\theta_1^{(V)})$  has multiple elements, then  $(j_1, r_1)$  is randomly drawn with probability

$$\mathbb{P}(\mathcal{P}_1 \in T^*(\Theta) | \Theta^{(V)} = \theta^{(V)}) = \Phi_{\theta_1^{(V)}, (j_1, r_1)}(\mathbf{0}), \quad (\text{C.1})$$

where  $\Phi_{\theta_1^{(V)}, (j_1, r_1)}$  is the cdf of the limit law defined in Lemma 4 for  $\mathcal{C}_1 = \mathcal{C}_1^*(\theta_1^{(V)})$ . This definition is extended for the second split in Definition 2.

Recall that the randomness in a tree can be decomposed as  $\Theta = (\Theta^{(S)}, \Theta^{(V)})$ , where  $\Theta^{(S)}$  corresponds to the subsampling and  $\Theta^{(V)}$  is related to the variable selection.  $\Theta^{(V)}$  takes values in the finite set  $\Omega^{(V)} = \{1, \dots, p\}^{3 \times \text{mtry}}$ .

**Lemma 5.** If Assumptions (A1)-(A3) are satisfied, then for all  $\theta^{(V)} \in \Omega^{(V)}$ , we have

$$\lim_{n \rightarrow \infty} \mathbb{P}(\mathcal{P}_1 \in T(\Theta, \mathcal{D}_n) | \Theta^{(V)} = \theta^{(V)}) = \mathbb{P}(\mathcal{P}_1 \in T^*(\Theta) | \Theta^{(V)} = \theta^{(V)}).$$

We are now equipped to prove Theorem 1 in the case of one single split. Recall that

$$\mathbb{E}[\hat{p}_{M_n, n}(\mathcal{P}_1)] = \mathbb{P}(\mathcal{P}_1 \in T(\Theta, \mathcal{D}_n)). \quad (\text{C.2})$$

Since  $\Theta^{(V)}$  takes values in the finite set  $\Omega^{(V)}$ , according to Lemma 5, we have

$$\begin{aligned} & \lim_{n \rightarrow \infty} \mathbb{P}(\mathcal{P}_1 \in T(\Theta, \mathcal{D}_n)) \\ &= \lim_{n \rightarrow \infty} \sum_{\theta^{(V)} \in \Omega^{(V)}} \mathbb{P}(\mathcal{P}_1 \in T(\Theta, \mathcal{D}_n) | \Theta^{(V)} = \theta^{(V)}) \mathbb{P}_{\Theta^{(V)}}(\Theta^{(V)} = \theta^{(V)}) \\ &= \sum_{\theta^{(V)} \in \Omega^{(V)}} \mathbb{P}(\mathcal{P}_1 \in T^*(\Theta) | \Theta^{(V)} = \theta^{(V)}) \mathbb{P}_{\Theta^{(V)}}(\Theta^{(V)} = \theta^{(V)}) \\ &= \mathbb{P}(\mathcal{P}_1 \in T^*(\Theta)). \end{aligned}$$

Therefore,

$$\lim_{n \rightarrow \infty} \mathbb{E}[\hat{p}_{M_n, n}(\mathcal{P}_1)] = p^*(\mathcal{P}_1).$$

To finish the proof, we just have to show that  $\lim_{n \rightarrow \infty} \mathbb{V}[\hat{p}_{M_n, n}(\mathcal{P}_1)] = 0$ .

The law of total variance gives

$$\begin{aligned} \mathbb{V}[\hat{p}_{M_n, n}(\mathcal{P}_1)] &= \mathbb{E}[\mathbb{V}[\hat{p}_{M_n, n}(\mathcal{P}_1)|\mathcal{D}_n]] + \mathbb{V}[\mathbb{E}[\hat{p}_{M_n, n}(\mathcal{P}_1)|\mathcal{D}_n]] \\ &= \mathbb{E}\left[\mathbb{V}\left[\frac{1}{M_n} \sum_{\ell=1}^{M_n} \mathbb{1}_{\mathcal{P}_1 \in T(\Theta_\ell, \mathcal{D}_n)} | \mathcal{D}_n\right]\right] + \mathbb{V}[p_n(\mathcal{P}_1)] \\ &= \frac{1}{M_n} \mathbb{E}[\mathbb{V}[\mathbb{1}_{\mathcal{P}_1 \in T(\Theta_1, \mathcal{D}_n)} | \mathcal{D}_n]] + \mathbb{V}[p_n(\mathcal{P}_1)] \\ &= \frac{1}{M_n} \mathbb{E}[p_n(\mathcal{P}_1) - p_n(\mathcal{P}_1)^2] + \mathbb{V}[p_n(\mathcal{P}_1)], \\ &= \frac{1}{M_n} \mathbb{E}[p_n(\mathcal{P}_1)](1 - \mathbb{E}[p_n(\mathcal{P}_1)]) + \left(1 - \frac{1}{M_n}\right) \mathbb{V}[p_n(\mathcal{P}_1)]. \end{aligned} \quad (\text{C.3})$$

Following the approach of [Mentch and Hooker \(2016\)](#),  $p_n(\mathcal{P}_1)$  is a complete infinite order U-statistic with the kernel  $\mathbb{E}[\mathbb{1}_{\mathcal{P}_1 \in T(\Theta, \mathcal{D}_n)} | \Theta^{(S)}, \mathcal{D}_n]$ . From [Hoeffding \(1948\)](#),

$$\mathbb{V}[p_n(\mathcal{P}_1)] \leq \frac{a_n}{n} \xi_{a_n, a_n},$$

where  $\xi_{a_n, a_n} = \mathbb{V}[\mathbb{E}[\mathbb{1}_{\mathcal{P}_1 \in T(\Theta, \mathcal{D}_n)} | \Theta^{(S)}, \mathcal{D}_n] | \Theta^{(S)}]$ . Since  $\xi_{a_n, a_n}$  is bounded and  $\lim_{n \rightarrow \infty} \frac{a_n}{n} = 0$ ,

$$\lim_{n \rightarrow \infty} \mathbb{V}[p_n(\mathcal{P}_1)] = 0.$$

Using equality (C.3), since  $p_n(\mathcal{P}_1)$  is bounded and  $\lim_{n \rightarrow \infty} M_n = \infty$ ,

$$\lim_{n \rightarrow \infty} \mathbb{V}[\hat{p}_{M_n, n}(\mathcal{P}_1)] = 0.$$

Finally,

$$\begin{aligned} &\lim_{n \rightarrow \infty} \mathbb{E}[(\hat{p}_{M_n, n}(\mathcal{P}_1) - p^*(\mathcal{P}_1))^2] \\ &= \lim_{n \rightarrow \infty} \mathbb{V}[\hat{p}_{M_n, n}(\mathcal{P}_1)] + (\mathbb{E}[\hat{p}_{M_n, n}(\mathcal{P}_1)] - p^*(\mathcal{P}_1))^2 = 0, \end{aligned}$$

which concludes the proof.

## C.2 Proof of Theorem 1 for a path of two split

The proof of Theorem 1 is extended for a path of two splits. We consider  $\mathcal{P}_1 = \{(j_1, r_1, s_1)\}$  a path of one split and  $\mathcal{P}_2 = \{(j_k, r_k, s_k), k = 1, 2\}$  a path of two splits, where  $j_1, j_2 \in \{1, \dots, p\}$ ,  $r_1, r_2 \in \{1, \dots, q-1\}$  and  $s_1, s_2 \in \{L, R\}$ . We assume assumptions (A1)-(A3) are satisfied.

The path  $\mathcal{P}_2 = \{(j_1, r_1, s_1), (j_2, r_2, s_2)\}$  can occur in trees where the split at the root node is  $(j_1, r_1)$  and the split of one of the child node is  $(j_2, r_2)$ , and in trees where the splits are

made in the reversed order,  $(j_2, r_2)$  at the root node and  $(j_1, r_1)$  at one of the child node. Since these two events are disjoint,  $\mathbb{P}(\mathcal{P}_2 \in T(\Theta, \mathcal{D}_n))$  is the sum of the probability of these two events. Without loss of generality, we will consider in the entire proof that the split at the root node is  $(j_1, r_1)$ . Lemmas 6 - 9 below extend Lemmas 2 - 5 to the case where the tree path contains two splits.

We need to introduce additional notations, first, the theoretical hyperrectangle based on a path  $\mathcal{P}$  by

$$H^*(\mathcal{P}) = \left\{ \mathbf{x} \in \mathbb{R}^p : \begin{cases} x^{(j_k)} < q_{r_k}^{*(j_k)} & \text{if } s_k = L \\ x^{(j_k)} \geq q_{r_k}^{*(j_k)} & \text{if } s_k = R \end{cases}, k \in 1, \dots, d \right\},$$

with  $d \in \{1, 2\}$ , the empirical counterpart of  $\hat{H}_n(\mathcal{P})$  defined in (2.3). Furthermore, since from assumption (A3),  $\mathbf{X}$  has a strictly positive density, then for  $j \in \{1, \dots, p\} \setminus j_1$ , and  $r \in \{1, \dots, q-1\}$ ,  $\mathbb{P}(\mathbf{X} \in H^*(\mathcal{P}_1), X^{(j)} < q_r^{*(j)}) > 0$  and  $\mathbb{P}(\mathbf{X} \in H^*(\mathcal{P}_1), X^{(j)} \geq q_r^{*(j)}) > 0$ . When  $j = j_1$ , the second cut is performed along the same direction as the first one. In that case, depending on the side  $s_1$  of the first cut and the cut positions  $r_1$  and  $r$ , one of the two child node can be empty with probability one. For example, the hyperrectangle associated to the path  $\{(1, 2, L), (1, 3, R)\}$  is empty. In SIRUS, such splits are not considered to find the best cut for a node at the second level of the tree. Thus we define  $\mathcal{C}_{\mathcal{P}_1}$  the set of possible splits for the second cut

$$\begin{aligned} \mathcal{C}_{\mathcal{P}_1} = & \{(j, r), j \in \{1, \dots, p\} \setminus j_1, r \in \{1, \dots, q-1\}\} \\ & \cup \{(j_1, r), \text{ s.t. } r < r_1 \text{ if } s_1 = L, \text{ and } r > r_1 \text{ if } s_1 = R\}, \end{aligned}$$

and  $\mathcal{C}_{\mathcal{P}_1}(\theta_2^{(V)}) = \{(j, r) \in \mathcal{C}_{\mathcal{P}_1} \text{ s.t. } j \in \theta_2^{(V)}\}$  when the split directions are restricted to  $\theta_2^{(V)} \subset \{1, \dots, p\}$ .

**Lemma 6.** *If Assumptions (A1) and (A3) are satisfied, then for all  $(j, r) \in \mathcal{C}_{\mathcal{P}_1}$ , we have*

$$\lim_{n \rightarrow \infty} \sqrt{a_n} (L_{a_n}(\hat{H}_n(\mathcal{P}_1), \hat{q}_{n,r}^{(j)}) - L_{a_n}(H^*(\mathcal{P}_1), q_r^{*(j)})) = 0 \quad \text{in probability.}$$

**Lemma 7.** *If Assumptions (A1) and (A3) are satisfied, then for all  $(j, r) \in \mathcal{C}_{\mathcal{P}_1}$ , we have*

$$\lim_{n \rightarrow \infty} L_{a_n}(\hat{H}_n(\mathcal{P}_1), \hat{q}_{n,r}^{(j)}) = L^*(H^*(\mathcal{P}_1), q_r^{*(j)}) \quad \text{in probability.}$$

**Lemma 8.** *Consider that Assumptions (A1) and (A3) are satisfied. Let  $\mathcal{C}_1 \subset \{1, \dots, p\} \times \{1, \dots, q-1\}$  and  $\mathcal{C}_2 \subset \mathcal{C}_{\mathcal{P}_1}$  be two sets of splits of cardinality  $c_1 \geq 1$  and  $c_2 \geq 2$ , such that the theoretical CART-splitting criterion is constant for all splits in  $\mathcal{C}_1$  on one hand, and in  $\mathcal{C}_2$  on the other hand, i.e.,*

$$\forall l \in \{1, 2\}, \quad \forall (j, r) \in \mathcal{C}_l, \quad L^*(H_l, q_r^{*(j)}) \stackrel{\text{def}}{=} L_{\mathcal{C}_l}^*,$$

where  $H_1 = \mathbb{R}^p$  and  $H_2 = H^*(\mathcal{P}_1)$ . Let  $(j_1, r_1) \in \mathcal{C}_1$ ,  $(j_2, r_2) \in \mathcal{C}_2$ , and let  $\mathbf{L}_{n, \mathcal{P}_2}^{(\mathcal{C}_1, \mathcal{C}_2)}$  a the random vector where each component is the difference between the empirical CART-splitting criterion for the splits  $(j, r) \in \mathcal{C}_1 \setminus (j_1, r_1)$  and  $(j_1, r_1)$  for the first  $c_1 - 1$  components, and for the splits  $(j, r) \in \mathcal{C}_2 \setminus (j_2, r_2)$  and  $(j_2, r_2)$  for the remaining  $c_2 - 1$  components, that is

$$\mathbf{L}_{n, \mathcal{P}_2}^{(\mathcal{C}_1, \mathcal{C}_2)} = \begin{pmatrix} [L_{a_n}(\mathbb{R}^p, \hat{q}_{n,r}^{(j)}) - L_{a_n}(\mathbb{R}^p, \hat{q}_{n,r_1}^{(j_1)})]_{(j,r) \in \mathcal{C}_1 \setminus (j_1, r_1)} \\ [L_{a_n}(\hat{H}_n(\mathcal{P}_1), \hat{q}_{n,r}^{(j)}) - L_{a_n}(\hat{H}_n(\mathcal{P}_1), \hat{q}_{n,r_2}^{(j_2)})]_{(j,r) \in \mathcal{C}_2 \setminus (j_2, r_2)} \end{pmatrix}.$$

(a) If  $L_{\mathcal{C}_1}^* > 0$  and  $L_{\mathcal{C}_2}^* > 0$ , then we have

$$\sqrt{a_n} \mathbf{L}_{n, \mathcal{P}_2}^{(\mathcal{C}_1, \mathcal{C}_2)} \xrightarrow[n \rightarrow \infty]{\mathcal{D}} \mathcal{N}(0, \Sigma)$$

where for  $l, l' \in \{1, 2\}$ , for all  $(j, r) \in \mathcal{C}_l \setminus (j_l, r_l)$ ,  $(j', r') \in \mathcal{C}_{l'} \setminus (j_{l'}, r_{l'})$ , each element of the covariance matrix  $\Sigma$  is defined by  $\Sigma_{(j,r,l),(j',r',l')} = \text{Cov}[Z_{j,r,l}, Z_{j',r',l'}]$ , with

$$\begin{aligned} Z_{j,r,l} = & \frac{1}{\mathbb{P}(\mathbf{X} \in H_l)} (Y - \mu_{L,r,l}^{(j)} \mathbb{1}_{X^{(j)} < q_{r_l}^{*(j)}} - \mu_{R,r,l}^{(j)} \mathbb{1}_{X^{(j)} \geq q_{r_l}^{*(j)}})^2 \mathbb{1}_{\mathbf{X} \in H_l} \\ & - \frac{1}{\mathbb{P}(\mathbf{X} \in H_l)} (Y - \mu_{L,r}^{(j)} \mathbb{1}_{X^{(j)} < q_r^{*(j)}} - \mu_{R,r}^{(j)} \mathbb{1}_{X^{(j)} \geq q_r^{*(j)}})^2 \mathbb{1}_{\mathbf{X} \in H_l}, \end{aligned}$$

$\mu_{L,r}^{(j)} = \mathbb{E}[Y | X^{(j)} < q_r^{*(j)}, \mathbf{X} \in H_l]$ ,  $\mu_{R,r}^{(j)} = \mathbb{E}[Y | X^{(j)} \geq q_r^{*(j)}, \mathbf{X} \in H_l]$ . Besides, for all  $l \in \{1, 2\}$  and for all  $(j, r) \in \mathcal{C}_l$ ,  $\mathbb{V}[Z_{j,r,l}] > 0$ .

(b) If  $L_{\mathcal{C}_1}^* = L_{\mathcal{C}_2}^* = 0$ , then we have

$$a_n \mathbf{L}_{n, \mathcal{P}_2}^{(\mathcal{C}_1, \mathcal{C}_2)} \xrightarrow[n \rightarrow \infty]{\mathcal{D}} h_{\mathcal{P}_2}(\mathbf{V}),$$

where  $\mathbf{V}$  is a gaussian vector of covariance matrix  $\text{Cov}[\mathbf{Z}]$ . If  $\mathcal{C}_1$  and  $\mathcal{C}_2$  are explicitly written  $\mathcal{C}_1 = \{(j_k, r_k)\}_{k \in J_1}$ , and  $\mathcal{C}_2 = \{(j_k, r_k)\}_{k \in J_2}$ , with  $J_1 = \{1, \dots, c_1 + 1\} \setminus 2$  and  $J_2 = \{2\} \cup \{c_1 + 2, \dots, c_1 + c_2\}$ ,  $\mathbf{Z}$  is defined, for  $l \in \{1, 2\}$  and  $k \in J_l$ , by

$$\begin{aligned} Z_{2k-1} &= \frac{1}{\sqrt{p_{L,k} \mathbb{P}(\mathbf{X} \in H_l)}} (Y - \mathbb{E}[Y | \mathbf{X} \in H_l]) \mathbb{1}_{X^{(j_k)} < q_{r_k}^{*(j_k)}} \mathbb{1}_{\mathbf{X} \in H_l} \\ Z_{2k} &= \frac{1}{\sqrt{p_{R,k} \mathbb{P}(\mathbf{X} \in H_l)}} (Y - \mathbb{E}[Y | \mathbf{X} \in H_l]) \mathbb{1}_{X^{(j_k)} \geq q_{r_k}^{*(j_k)}} \mathbb{1}_{\mathbf{X} \in H_l}, \end{aligned}$$

where  $p_{L,k} = \mathbb{P}(X^{(j_k)} < q_{r_k}^{*(j_k)}, \mathbf{X} \in H_l)$ ,  $p_{R,k} = \mathbb{P}(X^{(j_k)} \geq q_{r_k}^{*(j_k)}, \mathbf{X} \in H_l)$ , and  $h_{\mathcal{P}_2}$  is a multivariate quadratic form defined as

$$h_{\mathcal{P}_2} : \begin{pmatrix} x_1 \\ \vdots \\ x_{2(c_1+c_2)} \end{pmatrix} \rightarrow \begin{pmatrix} x_5^2 + x_6^2 - x_1^2 - x_2^2 \\ \vdots \\ x_{2c_1+1}^2 + x_{2c_1+2}^2 - x_1^2 - x_2^2 \\ x_{2c_1+3}^2 + x_{2c_1+4}^2 - x_3^2 - x_4^2 \\ \vdots \\ x_{2(c_1+c_2)-1}^2 + x_{2(c_1+c_2)}^2 - x_3^2 - x_4^2 \end{pmatrix}.$$

Besides, the variance of each component of  $h_{\mathcal{P}_2}(\mathbf{V})$  is strictly positive.

(c) If  $L_{\mathcal{C}_1}^* > 0$  and  $L_{\mathcal{C}_2}^* = 0$ , then we have

$$a_n \mathbf{L}_{n, \mathcal{P}_2}^{(\mathcal{C}_1, \mathcal{C}_2)} \xrightarrow[n \rightarrow \infty]{\mathcal{D}} h'_{\mathcal{P}_2}(\mathbf{V}),$$

where  $\mathbf{V}$  is a gaussian vector of covariance matrix  $\text{Cov}[\mathbf{Z}]$ , and  $\mathbf{Z}$  is defined as, for  $k \in J_1$ ,

$$\begin{aligned} Z_{2k-1} &= \left( Y - \mathbb{E}[Y|X^{(j_k)} < q_{r_k}^{*(j_k)}] \right)^2 \mathbb{1}_{X^{(j_k)} < q_{r_k}^{*(j_k)}} \\ Z_{2k} &= \left( Y - \mathbb{E}[Y|X^{(j_k)} \geq q_{r_k}^{*(j_k)}] \right)^2 \mathbb{1}_{X^{(j_k)} \geq q_{r_k}^{*(j_k)}}, \end{aligned}$$

for  $k \in J_2$ ,

$$\begin{aligned} Z_{2k-1} &= \frac{Y - \mathbb{E}[Y|\mathbf{X} \in H^*(\mathcal{P}_1)]}{\sqrt{p_{L,k} \mathbb{P}(\mathbf{X} \in H^*(\mathcal{P}_1))}} \mathbb{1}_{X^{(j_k)} < q_{r_k}^{*(j_k)}, \mathbf{X} \in H^*(\mathcal{P}_1)} \\ Z_{2k} &= \frac{Y - \mathbb{E}[Y|\mathbf{X} \in H^*(\mathcal{P}_1)]}{\sqrt{p_{R,k} \mathbb{P}(\mathbf{X} \in H^*(\mathcal{P}_1))}} \mathbb{1}_{X^{(j_k)} \geq q_{r_k}^{*(j_k)}, \mathbf{X} \in H^*(\mathcal{P}_1)}, \end{aligned}$$

and  $h'_{\mathcal{P}_2}$  is a multivariate quadratic form defined as

$$h'_{\mathcal{P}_2} : \begin{pmatrix} x_1 \\ \vdots \\ x_{2(c_1+c_2)} \end{pmatrix} \rightarrow \begin{pmatrix} x_1 + x_2 - x_5 - x_6 \\ \vdots \\ x_1 + x_2 - x_{2c_1+1} - x_{2c_1+2} \\ x_{2c_1+3}^2 + x_{2c_1+4}^2 - x_3^2 - x_4^2 \\ \vdots \\ x_{2(c_1+c_2)-1}^2 + x_{2(c_1+c_2)}^2 - x_3^2 - x_4^2 \end{pmatrix}.$$

Besides, the variance of each component of  $h'_{\mathcal{P}_2}(\mathbf{V})$  is strictly positive.

(d)  $L_{\mathcal{C}_1}^* = 0$  and  $L_{\mathcal{C}_2}^* > 0$ . Symmetric to case (c).

**Definition 2** (Theoretical splitting procedure at children nodes). Let  $\theta^{(V)} = (\theta_1^{(V)}, \theta_2^{(V)}, \cdot) \in \Omega^{(V)}$  be the sets of eligible variables to split the nodes of a theoretical random tree. The set of best theoretical cuts at the left children node along the variables in  $\theta_2^{(V)}$  is defined as

$$\mathcal{C}_2^*(\theta_2^{(V)}) = \underset{(j,r) \in \mathcal{C}_{\mathcal{P}_1}(\theta_2^{(V)})}{\operatorname{argmax}} L^*(H^*(\mathcal{P}_1), q_r^{*(j)}).$$

If  $\mathcal{C}_2^*(\theta_2^{(V)})$  has multiple elements, then  $(j_2, r_2)$  is randomly drawn with probability

$$\mathbb{P}(\mathcal{P}_2 \in T^*(\Theta) | \Theta^{(V)} = \theta^{(V)}) = \frac{\Phi_{\mathcal{P}_1, \theta^{(V)}, (j_2, r_2)}(\boldsymbol{\theta})}{\mathbb{P}(\mathcal{P}_1 \in T^*(\Theta) | \Theta^{(V)} = \theta^{(V)})}, \quad (\text{C.4})$$

where  $\mathbb{P}(\mathcal{P}_1 \in T^*(\Theta) | \Theta^{(V)} = \theta^{(V)})$  is defined from Definition 1, and  $\Phi_{\mathcal{P}_1, \theta^{(V)}, (j_2, r_2)}$  is the cdf of the limit law defined in Lemma 8 for  $\mathcal{C}_1 = \mathcal{C}_1^*(\theta_1^{(V)})$  and  $\mathcal{C}_2 = \mathcal{C}_2^*(\theta_2^{(V)})$ .

**Lemma 9.** If Assumptions (A1)-(A3) are satisfied, then for all  $\theta^{(V)} \in \Omega^{(V)}$ , we have

$$\lim_{n \rightarrow \infty} \mathbb{P}(\mathcal{P}_2 \in T(\Theta, \mathcal{D}_n) | \Theta^{(V)} = \theta^{(V)}) = \mathbb{P}(\mathcal{P}_2 \in T^*(\Theta) | \Theta^{(V)} = \theta^{(V)})$$

Finally, the proof of Theorem 1 in the two-splits scenario is the same as in the single-split scenario.

### C.3 Proofs of intermediate lemmas

*Proof of Lemma 1.* Set  $j \in \{1, \dots, p\}$ , and  $r \in \{1, \dots, q-1\}$ . We define the marginal cumulative distribution function  $F^{(j)}$  of  $X^{(j)}$ ,  $F^{(j)}(x) = \mathbb{P}(X^{(j)} < x)$ , and  $F_n^{(j)}$  the empirical c.d.f.

$$F_n^{(j)}(x) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{X_i^{(j)} \leq x}.$$

We adapt an inequality from [Serfling \(2009\)](#) (section 2.3.2 page 75) to bound the following conditional probability for all  $\varepsilon > 0$

$$\begin{aligned} & \mathbb{P}(q_r^{*(j)} \leq X_1^{(j)} < \hat{q}_{n,r}^{(j)} | X_1^{(j)} = q_r^{*(j)} + \varepsilon) \\ &= \mathbb{P}(q_r^{*(j)} + \varepsilon < \hat{q}_{n,r}^{(j)} | X_1^{(j)} = q_r^{*(j)} + \varepsilon) \\ &\leq \mathbb{P}(F_n^{(j)}(q_r^{*(j)} + \varepsilon) \leq F_n^{(j)}(\hat{q}_{n,r}^{(j)}) | X_1^{(j)} = q_r^{*(j)} + \varepsilon) \\ &\leq \mathbb{P}\left(1 + \sum_{i=2}^n \mathbb{1}_{X_i^{(j)} \leq q_r^{*(j)} + \varepsilon} \leq \left\lceil \frac{n \cdot r}{q} \right\rceil\right) \\ &\leq \mathbb{P}\left(\sum_{i=2}^n \mathbb{1}_{X_i^{(j)} \leq q_r^{*(j)} + \varepsilon} - (n-1)F^{(j)}(q_r^{*(j)} + \varepsilon)\right) \end{aligned} \tag{C.5}$$

$$\leq \left\lceil \frac{n \cdot r}{q} \right\rceil - 1 - (n-1)F^{(j)}(q_r^{*(j)} + \varepsilon) \tag{C.6}$$

Since  $f$  is continuous and strictly positive, there exists three constants  $c_1, c_2, \eta > 0$  such that for all  $x \in [q_r^{*(j)}, q_r^{*(j)} + \eta]$ ,  $c_1 \leq f^{(j)}(x) \leq c_2$ . Thus, for all  $\varepsilon < \eta$ , we have

$$F^{(j)}(q_r^{*(j)} + \varepsilon) - F^{(j)}(q_r^{*(j)}) = \int_{q_r^{*(j)}}^{q_r^{*(j)} + \varepsilon} f^{(j)}(x) dx,$$

which leads to

$$c_1 \varepsilon \leq F^{(j)}(q_r^{*(j)} + \varepsilon) - F^{(j)}(q_r^{*(j)}) \leq c_2 \varepsilon.$$

Consequently,

$$\begin{aligned} & \left\lceil \frac{n \cdot r}{q} \right\rceil - 1 - (n-1)F^{(j)}(q_r^{*(j)} + \varepsilon) \\ &\leq \left\lceil \frac{n \cdot r}{q} \right\rceil - 1 - (n-1)(c_1 \varepsilon + F^{(j)}(q_r^{*(j)})) \\ &\leq \left\lceil \frac{n \cdot r}{q} \right\rceil - 1 - (n-1)c_1 \varepsilon - \frac{(n-1) \cdot r}{q} \\ &\leq 1 - (n-1)c_1 \varepsilon. \end{aligned}$$

For  $n > 1 + \frac{1}{c_1 \varepsilon}$ , we can apply Hoeffding inequality to [C.6](#),

$$\begin{aligned}
& \mathbb{P}(q_r^{*(j)} \leq X_1^{(j)} < \hat{q}_{n,r}^{(j)} | X_1^{(j)} = q_r^{*(j)} + \varepsilon) \\
& \leq \mathbb{P}\left(\sum_{i=2}^n \mathbf{1}_{X_i^{(j)} \leq q_r^{*(j)} + \varepsilon} - (n-1)F^{(j)}(q_r^{*(j)} + \varepsilon) \leq 1 - (n-1)c_1\varepsilon\right) \\
& \leq e^{-\frac{2}{n}(1-(n-1)c_1\varepsilon)^2} \\
& \leq Ce^{-2nc_1^2\varepsilon^2},
\end{aligned} \tag{C.7}$$

where  $C = e^{2c_1\eta(1+2c_1\eta)}$ . By definition, we have

$$\begin{aligned}
\mathbb{P}(q_r^{*(j)} \leq X_1^{(j)} < \hat{q}_{n,r}^{(j)}) &= \int_{]0, \infty[} \mathbb{P}(q_r^{*(j)} \leq X_1^{(j)} < \hat{q}_{n,r}^{(j)} | X_1^{(j)} = q_r^{*(j)} + \varepsilon) \\
&\quad \times f^{(j)}(q_r^{*(j)} + \varepsilon) d\varepsilon.
\end{aligned}$$

To bound the previous integral, we break it down in three parts. Since  $f^{(j)}$  is bounded by  $c_2$  on  $[q_r^{*(j)}, q_r^{*(j)} + \eta]$ , for  $n > 1 + \frac{1}{c_1\eta}$  we use inequality [C.7](#) to get

$$\begin{aligned}
\mathbb{P}(q_r^{*(j)} \leq X_1^{(j)} < \hat{q}_{n,r}^{(j)}) &\leq \int_{]0, \frac{1}{(n-1)c_1}]} c_2 d\varepsilon \\
&\quad + \int_{] \frac{1}{(n-1)c_1}, \eta]} c_2 Ce^{-2nc_1^2\varepsilon^2} d\varepsilon \\
&\quad + \int_{[\eta, \infty[} Ce^{-2nc_1^2\eta^2} f^{(j)}(q_r^{*(j)} + \varepsilon) d\varepsilon.
\end{aligned}$$

In the second integral, we introduce the following change of variable  $u = \sqrt{2n}c_1\varepsilon$

$$\begin{aligned}
\int_{] \frac{1}{(n-1)c_1}, \eta]} c_2 Ce^{-2nc_1^2\varepsilon^2} d\varepsilon &= \frac{c_2 C}{c_1 \sqrt{2n}} \int_{] \frac{\sqrt{2n}}{(n-1)}, \sqrt{2n}c_1\eta]} e^{-u^2} du \\
&\leq \frac{c_2 C}{c_1 \sqrt{2n}} \int_{]0, \infty[} e^{-u^2} du \leq \frac{\sqrt{\pi}c_2 C}{2c_1 \sqrt{2n}},
\end{aligned}$$

and therefore we can write

$$\sqrt{a_n} \mathbb{P}(q_r^{*(j)} \leq X_1^{(j)} < \hat{q}_{n,r}^{(j)}) \leq \frac{c_2 \sqrt{a_n}}{(n-1)c_1} + \frac{\sqrt{\pi a_n} c_2 C}{2c_1 \sqrt{2n}} + C \sqrt{a_n} e^{-2nc_1^2\eta^2}$$

From Assumption (A1),  $\lim_{n \rightarrow \infty} \frac{a_n}{n} = 0$ , and then

$$\lim_{n \rightarrow \infty} \sqrt{a_n} \mathbb{P}(q_r^{*(j)} \leq X_1^{(j)} < \hat{q}_{n,r}^{(j)}) = 0.$$

The case  $\lim_{n \rightarrow \infty} \sqrt{a_n} \mathbb{P}(\hat{q}_{n,r}^{(j)} \leq X_1^{(j)} < q_r^{*(j)}) = 0$  is similar. □



### C.3.1 Case 1: $\mathcal{P}_1$

*Proof of Lemma 2.* Let  $j \in \{1, \dots, p\}$ ,  $r \in \{1, \dots, q-1\}$ , and  $H \subseteq \mathbb{R}^p$  such that  $\mathbb{P}(\mathbf{X} \in H, X^{(j)} < q_r^{\star(j)}) > 0$  and  $\mathbb{P}(\mathbf{X} \in H, X^{(j)} \geq q_r^{\star(j)}) > 0$ . Let

$$\Delta_{n,r}^{(j)} = \sqrt{a_n} (L_{a_n}(H, \hat{q}_{n,r}^{(j)}) - L_{a_n}(H, q_r^{\star(j)}))$$

that is

$$\begin{aligned} \Delta_{n,r}^{(j)} = & -\frac{\sqrt{a_n}}{N_n(H)} \left[ \sum_{i=1}^{a_n} (Y_i - \bar{Y}_{H_L} \mathbb{1}_{X_i^{(j)} < \hat{q}_{n,r}^{(j)}} - \bar{Y}_{H_R} \mathbb{1}_{X_i^{(j)} \geq \hat{q}_{n,r}^{(j)}})^2 \mathbb{1}_{\mathbf{X}_i \in H} \right. \\ & \left. - \sum_{i=1}^{a_n} (Y_i - \bar{Y}_{H_L^*} \mathbb{1}_{X_i^{(j)} < q_r^{\star(j)}} - \bar{Y}_{H_R^*} \mathbb{1}_{X_i^{(j)} \geq q_r^{\star(j)}})^2 \mathbb{1}_{\mathbf{X}_i \in H} \right] \end{aligned}$$

where, for a generic hyperrectangle  $H$ , we define  $N_n(H) = \sum_{i=1}^{a_n} \mathbb{1}_{\mathbf{X}_i \in H}$ , and

$$H_L = \{\mathbf{x} \in H : x^{(j)} < \hat{q}_{n,r}^{(j)}\} \quad \text{and} \quad \bar{Y}_{H_L} = \frac{1}{N_n(H_L)} \sum_{i=1}^{a_n} Y_i \mathbb{1}_{X_i^{(j)} < \hat{q}_{n,r}^{(j)}} \mathbb{1}_{\mathbf{X}_i \in H},$$

with the convention  $\bar{Y}_{H_L} = 0$  if  $H_L$  is empty. The theoretical quantities  $H_L^*$  and  $\bar{Y}_{H_L^*}$  are defined similarly by replacing the empirical quantile by its population version. We define symmetrically  $H_R$ ,  $H_R^*$ ,  $\bar{Y}_{H_R}$ ,  $\bar{Y}_{H_R^*}$ .

Simple calculations show that

$$\begin{aligned} \Delta_{n,r}^{(j)} = & \frac{\sqrt{a_n}}{N_n(H)} (\bar{Y}_{H_L}^2 N_n(H_L) - \bar{Y}_{H_L^*}^2 N_n(H_L^*)) \\ & + \frac{\sqrt{a_n}}{N_n(H)} (\bar{Y}_{H_R}^2 N_n(H_R) - \bar{Y}_{H_R^*}^2 N_n(H_R^*)) \end{aligned} \quad (\text{C.8})$$

The first term in equation (C.8) can be rewritten as

$$\begin{aligned} & \frac{\sqrt{a_n}}{N_n(H)} (\bar{Y}_{H_L}^2 N_n(H_L) - \bar{Y}_{H_L^*}^2 N_n(H_L^*)) \\ & = \frac{\sqrt{a_n}}{N_n(H) N_n(H_L) N_n(H_L^*)} \sum_{i,k,l=1}^{a_n} Y_i Y_k \mathbb{1}_{\mathbf{X}_i \in H, \mathbf{X}_k \in H} \\ & \quad \times (\mathbb{1}_{X_l^{(j)} < q_r^{\star(j)}} \mathbb{1}_{X_i^{(j)} < \hat{q}_{n,r}^{(j)}} \mathbb{1}_{X_k^{(j)} < \hat{q}_{n,r}^{(j)}} - \mathbb{1}_{X_l^{(j)} < \hat{q}_{n,r}^{(j)}} \mathbb{1}_{X_i^{(j)} < q_r^{\star(j)}} \mathbb{1}_{X_k^{(j)} < q_r^{\star(j)}}). \end{aligned}$$

Since  $Y_i \in \{0, 1\}$ , we have the following bound

$$\begin{aligned} & \frac{\sqrt{a_n}}{N_n(H)} |\bar{Y}_{H_L}^2 N_n(H_L) - \bar{Y}_{H_L^*}^2 N_n(H_L^*)| \\ & \leq \frac{\sqrt{a_n}}{N_n(H) N_n(H_L) N_n(H_L^*)} \sum_{i,k,l=1}^{a_n} |\mathbb{1}_{X_l^{(j)} < q_r^{\star(j)}} \mathbb{1}_{X_i^{(j)} < \hat{q}_{n,r}^{(j)}} \mathbb{1}_{X_k^{(j)} < \hat{q}_{n,r}^{(j)}} \\ & \quad - \mathbb{1}_{X_l^{(j)} < \hat{q}_{n,r}^{(j)}} \mathbb{1}_{X_i^{(j)} < q_r^{\star(j)}} \mathbb{1}_{X_k^{(j)} < q_r^{\star(j)}}|, \end{aligned}$$

and finally

$$\frac{\sqrt{a_n}}{N_n(H)} |\bar{Y}_{H_L}^2 N_n(H_L) - \bar{Y}_{H_L^*}^2 N_n(H_L^*)| \leq \frac{a_n^3}{N_n(H)N_n(H_L)N_n(H_L^*)} W_{n,r}^{(j)}, \quad (\text{C.9})$$

where

$$W_{n,r}^{(j)} = \frac{\sqrt{a_n}}{a_n^3} \sum_{i,k,l=1}^{a_n} \left| \mathbb{1}_{X_l^{(j)} < q_r^{*(j)}} \mathbb{1}_{X_i^{(j)} < \hat{q}_{n,r}^{(j)}} \mathbb{1}_{X_k^{(j)} < \hat{q}_{n,r}^{(j)}} - \mathbb{1}_{X_l^{(j)} < \hat{q}_{n,r}^{(j)}} \mathbb{1}_{X_i^{(j)} < q_r^{*(j)}} \mathbb{1}_{X_k^{(j)} < q_r^{*(j)}} \right|. \quad (\text{C.10})$$

A close inspection of the terms inside the sum of (C.10) reveals that

$$\begin{aligned} \mathbb{E}[W_{n,r}^{(j)}] &\leq \frac{\sqrt{a_n}}{a_n^3} \sum_{i,k,l=1}^{a_n} \mathbb{P}(\hat{q}_{n,r}^{(j)} \leq X_i^{(j)} < q_r^{*(j)}) + \mathbb{P}(\hat{q}_{n,r}^{(j)} \leq X_k^{(j)} < q_r^{*(j)}) \\ &\quad + \mathbb{P}(q_r^{*(j)} \leq X_l^{(j)} < \hat{q}_{n,r}^{(j)}) + \mathbb{P}(q_r^{*(j)} \leq X_i^{(j)} < \hat{q}_{n,r}^{(j)}) \\ &\quad + \mathbb{P}(q_r^{*(j)} \leq X_k^{(j)} < \hat{q}_{n,r}^{(j)}) + \mathbb{P}(\hat{q}_{n,r}^{(j)} \leq X_l^{(j)} < q_r^{*(j)}) \\ &\leq 3\sqrt{a_n} \mathbb{P}(\hat{q}_{n,r}^{(j)} \leq X_1^{(j)} < q_r^{*(j)}) + 3\sqrt{a_n} \mathbb{P}(q_r^{*(j)} \leq X_1^{(j)} < \hat{q}_{n,r}^{(j)}), \end{aligned}$$

which tends to zero, according to Lemma 1. Thus, in probability,

$$\lim_{n \rightarrow \infty} W_{n,r}^{(j)} = 0. \quad (\text{C.11})$$

Regarding the remaining terms in inequality (C.9), by the law of large numbers, in probability,

$$\lim_{n \rightarrow \infty} \frac{N_n(H)}{a_n} = \mathbb{P}(\mathbf{X} \in H), \quad \lim_{n \rightarrow \infty} \frac{N_n(H_L^*)}{a_n} = \mathbb{P}(\mathbf{X} \in H_L^*). \quad (\text{C.12})$$

Additionally,

$$\begin{aligned} \mathbb{E}\left[\left|\frac{N_n(H_L)}{a_n} - \frac{N_n(H_L^*)}{a_n}\right|\right] &\leq \mathbb{E}\left[\frac{1}{a_n} \sum_{i=1}^{a_n} \mathbb{1}_{X_i^{(j)} \in H} \left| \mathbb{1}_{X_i^{(j)} \leq \hat{q}_{n,r}^{(j)}} - \mathbb{1}_{X_i^{(j)} \leq q_r^{*(j)}} \right|\right] \\ &\leq \mathbb{P}(\hat{q}_{n,r}^{(j)} \leq X_1^{(j)} < q_r^{*(j)}) + \mathbb{P}(q_r^{*(j)} \leq X_1^{(j)} < \hat{q}_{n,r}^{(j)}), \end{aligned}$$

which tends to zero, according to Lemma 1. Therefore, in probability,

$$\lim_{n \rightarrow \infty} \frac{N_n(H_L)}{a_n} - \frac{N_n(H_L^*)}{a_n} = 0. \quad (\text{C.13})$$

Since  $\mathbb{P}(\mathbf{X} \in H) > 0$  and  $\mathbb{P}(\mathbf{X} \in H_L^*) > 0$  by assumption, we can combine (C.11)-(C.13) to obtain, in probability,

$$\lim_{n \rightarrow \infty} \frac{a_n^3}{N_n(H)N_n(H_L)N_n(H_L^*)} = \frac{1}{\mathbb{P}(\mathbf{X} \in H)\mathbb{P}(\mathbf{X} \in H_L^*)^2}. \quad (\text{C.14})$$

Using (C.11) and (C.14) and inequality (C.9), we obtain, in probability,

$$\lim_{n \rightarrow \infty} \frac{\sqrt{a_n}}{N_n(H)} |\bar{Y}_{H_L}^2 N_n(H_L) - \bar{Y}_{H_L^*}^2 N_n(H_L^*)| = 0.$$

Similar results can be derived for the other term in equation (C.8), which allows us to conclude that, in probability,

$$\lim_{n \rightarrow \infty} \sqrt{a_n} (L_{a_n}(H, \hat{q}_{n,r}^{(j)}) - L_{a_n}(H, q_r^{*(j)})) = 0.$$

□

*Proof of Lemma 3.* Let  $j \in \{1, \dots, p\}$ ,  $r \in \{1, \dots, q-1\}$  and  $H \subseteq \mathbb{R}^p$  such that  $\mathbb{P}(\mathbf{X} \in H, X^{(j)} < q_r^{*(j)}) > 0$  and  $\mathbb{P}(\mathbf{X} \in H, X^{(j)} \geq q_r^{*(j)}) > 0$ .

$$L_{a_n}(H, \hat{q}_{n,r}^{(j)}) = L_{a_n}(H, q_r^{*(j)}) + (L_{a_n}(H, \hat{q}_{n,r}^{(j)}) - L_{a_n}(H, q_r^{*(j)}))$$

From the law of large number, in probability,

$$\lim_{n \rightarrow \infty} L_{a_n}(H, q_r^{*(j)}) = L^*(H, q_r^{*(j)}).$$

Thus, according to Lemma 2, in probability,

$$\lim_{n \rightarrow \infty} L_{a_n}(H, \hat{q}_{n,r}^{(j)}) = L^*(H, q_r^{*(j)}).$$

□

*Proof of Lemma 4.* We consider  $\mathcal{C}_1$ , a set of splits of cardinality  $c_1 \geq 2$  satisfying, for all  $(j, r) \in \mathcal{C}_1$ ,  $L^*(\mathbb{R}^p, q_r^{*(j)}) \stackrel{\text{def}}{=} L_{\mathcal{C}_1}^*$ . Fix  $(j_1, r_1) \in \mathcal{C}_1$ , we recall that

$$\mathbf{L}_{n, \mathcal{P}_1}^{(\mathcal{C}_1)} = \left( L_{a_n}(\mathbb{R}^p, \hat{q}_{n,r}^{(j)}) - L_{a_n}(\mathbb{R}^p, \hat{q}_{n,r_1}^{(j_1)}) \right)_{(j,r) \in \mathcal{C}_1 \setminus (j_1, r_1)}.$$

**Case (a):**  $L_{\mathcal{C}_1}^* > 0$  We first consider the following decomposition for  $(j, r) \in \mathcal{C}_1$ ,

$$\begin{aligned} L_{a_n}(\mathbb{R}^p, \hat{q}_{n,r}^{(j)}) &= L_{a_n}(\mathbb{R}^p, q_r^{*(j)}) + (L_{a_n}(\mathbb{R}^p, \hat{q}_{n,r}^{(j)}) - L_{a_n}(\mathbb{R}^p, q_r^{*(j)})) \\ &= \frac{1}{a_n} \sum_{i=1}^{a_n} (Y_i - \bar{Y})^2 - \frac{1}{a_n} \sum_{i=1}^{a_n} (Y_i - \bar{Y}_L^* \mathbb{1}_{X_i^{(j)} < q_r^{*(j)}} - \bar{Y}_R^* \mathbb{1}_{X_i^{(j)} \geq q_r^{*(j)}})^2 \\ &\quad + L_{a_n}(\mathbb{R}^p, \hat{q}_{n,r}^{(j)}) - L_{a_n}(\mathbb{R}^p, q_r^{*(j)}), \end{aligned}$$

where

$$N_{n,L}^* = \sum_{i=1}^{a_n} \mathbb{1}_{X_i^{(j)} < q_r^{*(j)}} \quad \text{and} \quad \bar{Y}_L^* = \frac{1}{N_{n,L}^*} \sum_{i=1}^{a_n} Y_i \mathbb{1}_{X_i^{(j)} < q_r^{*(j)}}$$

( $\bar{Y}_R^*$ ,  $N_{n,R}^*$  are defined symmetrically). Letting  $\mu_{L,r}^{(j)} = \mathbb{E}[Y | X^{(j)} < q_r^{*(j)}]$  (and  $\mu_{R,r}^{(j)}$  symmetrically), the first two terms of the last decomposition are standard variance estimates and we can write

$$L_{a_n}(\mathbb{R}^p, \hat{q}_{n,r}^{(j)}) = \frac{1}{a_n} \sum_{i=1}^{a_n} (Y_i - \bar{Y})^2 \tag{C.15}$$

$$- \frac{1}{a_n} \sum_{i=1}^{a_n} (Y_i - \mu_{L,r}^{(j)} \mathbb{1}_{X_i^{(j)} < q_r^{*(j)}} - \mu_{R,r}^{(j)} \mathbb{1}_{X_i^{(j)} \geq q_r^{*(j)}})^2 + R_{n,r}^{(j)}, \tag{C.16}$$

where

$$R_{n,L}^{(j)} = \frac{N_{n,L}^*}{a_n} (\bar{Y}_L^* - \mu_{L,r}^{(j)})^2 + \frac{N_{n,R}^*}{a_n} (\bar{Y}_R^* - \mu_{L,r}^{(j)})^2 + L_{a_n}(\mathbb{R}^p, \hat{q}_{n,r}^{(j)}) - L_{a_n}(\mathbb{R}^p, q_r^{*(j)}). \quad (\text{C.17})$$

Using the Central limit theorem, in probability,

$$\lim_{n \rightarrow \infty} \sqrt{a_n} \frac{N_{L,r}^*}{a_n} (\bar{Y}_{L,r}^* - \mu_{L,r}^{(j)})^2 = 0. \quad (\text{C.18})$$

The same result holds for the second term of (C.17), and using Lemma 2 for the third term of (C.17), we get that, in probability,

$$\lim_{n \rightarrow \infty} \sqrt{a_n} (L_{a_n}(\mathbb{R}^p, \hat{q}_{n,r}^{(j)}) - L_{a_n}(\mathbb{R}^p, q_r^{*(j)})) = 0.$$

Finally,

$$\lim_{n \rightarrow \infty} \sqrt{a_n} R_{n,r}^{(j)} = 0, \quad \text{in probability.}$$

Using Equation (C.16), each component of  $\mathbf{L}_{n,\mathcal{P}_1}^{(\mathcal{C}_1)}$  writes, with  $(j, r) \in \mathcal{C}_1 \setminus (j_1, r_1)$ ,

$$\begin{aligned} & L_{a_n}(\mathbb{R}^p, \hat{q}_{n,r}^{(j)}) - L_{a_n}(\mathbb{R}^p, \hat{q}_{n,r_1}^{(j_1)}) \\ &= \frac{1}{a_n} \sum_{i=1}^{a_n} (Y_i - \mu_{L,r_1}^{(j_1)} \mathbb{1}_{X_i^{(j_1)} < q_{r_1}^{*(j_1)}} - \mu_{R,r_1}^{(j_1)} \mathbb{1}_{X_i^{(j_1)} \geq q_{r_1}^{*(j_1)}})^2 \\ & \quad - (Y_i - \mu_{L,r}^{(j)} \mathbb{1}_{X_i^{(j)} < q_r^{*(j)}} - \mu_{R,r}^{(j)} \mathbb{1}_{X_i^{(j)} \geq q_r^{*(j)}})^2 \\ & \quad + R_{n,r}^{(j)} - R_{n,r_1}^{(j_1)} \end{aligned}$$

We can apply the multivariate Central limit theorem and Slutsky's theorem to obtain,

$$\sqrt{a_n} \mathbf{L}_{n,\mathcal{P}_1}^{(\mathcal{C}_1)} \xrightarrow[n \rightarrow \infty]{\mathcal{D}} \mathcal{N}(0, \Sigma)$$

where for all  $(j, r), (j', r') \in \mathcal{C}_1 \setminus (j_1, r_1)$ , each element of the covariance matrix  $\Sigma$  is defined by  $\Sigma_{(j,r),(j',r')} = \text{Cov}[Z_{j,r}, Z_{j',r'}]$ , with

$$\begin{aligned} Z_{j,r} &= (Y - \mu_{L,r_1}^{(j_1)} \mathbb{1}_{X^{(j_1)} < q_{r_1}^{*(j_1)}} - \mu_{R,r_1}^{(j_1)} \mathbb{1}_{X^{(j_1)} \geq q_{r_1}^{*(j_1)}})^2 \\ & \quad - (Y - \mu_{L,r}^{(j)} \mathbb{1}_{X^{(j)} < q_r^{*(j)}} - \mu_{R,r}^{(j)} \mathbb{1}_{X^{(j)} \geq q_r^{*(j)}})^2. \end{aligned}$$

Since  $L_{\mathcal{C}_1}^* > 0$ , we have for all  $(j, r) \in \mathcal{C}_1$ ,  $\mu_{L,r}^{(j)} \neq \mu_{R,r}^{(j)}$ . Besides, according to assumption (A3),  $\mathbf{X}$  has a strictly positive density. Consequently, the variance of  $Z_{j,r}$  is strictly positive. This concludes the first case.

**Case (b):**  $L_{\mathcal{C}_1}^* = 0$  Fix  $(j, r) \in \mathcal{C}_1$ . Since  $L^*(\mathbb{R}^p, q_r^{*(j)}) = 0$ , we have

$$\mathbb{E}[Y] = \mathbb{E}[Y|X^{(j)} < q_r^{*(j)}] = \mathbb{E}[Y|X^{(j)} \geq q_r^{*(j)}] \stackrel{\text{def}}{=} \mu.$$

Then, simple calculations show that  $L_{a_n}(\mathbb{R}^p, \hat{q}_{n,r}^{(j)})$  writes

$$L_{a_n}(\mathbb{R}^p, \hat{q}_{n,r}^{(j)}) = -(\bar{Y} - \mu)^2 + \underbrace{\frac{N_{n,L}}{a_n}(\bar{Y}_L - \mu)^2}_{\delta_L} + \underbrace{\frac{N_{n,R}}{a_n}(\bar{Y}_R - \mu)^2}_{\delta_R},$$

where

$$N_{n,L} = \sum_{i=1}^{a_n} \mathbb{1}_{X_i^{(j)} < \hat{q}_{n,r}^{(j)}} \quad \text{and} \quad \bar{Y}_L = \frac{1}{N_{n,L}} \sum_{i=1}^{a_n} Y_i \mathbb{1}_{X_i^{(j)} < \hat{q}_{n,r}^{(j)}}$$

( $N_{n,R}$ ,  $\bar{Y}_R$  are defined similarly for the other cell). Letting  $p_{L,r}^{(j)} = \mathbb{P}(X^{(j)} < q_r^{*(j)})$  and  $p_{R,r}^{(j)} = \mathbb{P}(X^{(j)} \geq q_r^{*(j)})$  with  $p_{L,r}^{(j)}, p_{R,r}^{(j)} > 0$ , we have

$$\begin{aligned} \delta_L &= \frac{N_{n,L}}{a_n} (\bar{Y}_L - \mu)^2 \\ &= \frac{N_{n,L}}{a_n} (\bar{Y}_L^* - \mu)^2 - 2 \frac{N_{n,L}}{a_n} (\bar{Y}_L^* - \bar{Y}_L) (\bar{Y}_L^* - \mu) + \frac{N_{n,L}}{a_n} (\bar{Y}_L^* - \bar{Y}_L)^2 \\ &= \frac{1}{p_{L,r}^{(j)}} \left( \frac{1}{a_n} \sum_{i=1}^{a_n} (Y_i - \mu) \mathbb{1}_{X_i^{(j)} < q_r^{*(j)}} \right)^2 + R_{L,r}^{(j)}, \end{aligned}$$

where

$$\begin{aligned} R_{L,r}^{(j)} &= \left( \frac{a_n N_{n,L}}{N_{n,L}^{*2}} - \frac{1}{p_{n,L}} \right) \left( \frac{1}{a_n} \sum_{i=1}^{a_n} (Y_i - \mu) \mathbb{1}_{X_i^{(j)} < q_r^{*(j)}} \right)^2 \\ &\quad - 2 \frac{N_{n,L}}{a_n} (\bar{Y}_L^* - \bar{Y}_L) (\bar{Y}_L^* - \mu) + \frac{N_{n,L}}{a_n} (\bar{Y}_L^* - \bar{Y}_L)^2 \end{aligned}$$

By the law of large numbers,  $\lim_{n \rightarrow \infty} \frac{N_{n,L}^*}{a_n} = p_{L,r}^{(j)}$  in probability. Using Equation (C.13) in the proof of Lemma 2, it comes that, in probability,  $\lim_{n \rightarrow \infty} \frac{N_{n,L}}{a_n} = p_{L,r}^{(j)}$ , and consequently  $\lim_{n \rightarrow \infty} \frac{a_n N_{n,L}}{N_{n,L}^{*2}} = \frac{1}{p_{L,r}^{(j)}}$ . Since  $\sqrt{a_n} \frac{1}{a_n} \sum_{i=1}^{a_n} (Y_i - \mu) \mathbb{1}_{X_i^{(j)} < q_r^{*(j)}}$  converges in distribution to a normal distribution by the Central limit theorem,

$$\lim_{n \rightarrow \infty} a_n \left( \frac{a_n N_{n,L}}{N_{n,L}^{*2}} - \frac{1}{p_{L,r}^{(j)}} \right) \left( \frac{1}{a_n} \sum_{i=1}^{a_n} (Y_i - \mu) \mathbb{1}_{X_i^{(j)} < q_r^{*(j)}} \right)^2 = 0, \quad \text{in probability.}$$

Furthermore, as for Equation (C.10) in the proof of Lemma 2,

$$\begin{aligned} &\sqrt{a_n} |\bar{Y}_L^* - \bar{Y}_L| \\ &\leq \underbrace{\frac{a_n^2}{N_{n,L} N_{n,L}^*} \frac{\sqrt{a_n}}{a_n^2} \sum_{i=1, l=1}^{a_n} Y_i \left| \mathbb{1}_{X_i^{(j)} < q_r^{*(j)}} \mathbb{1}_{X_l^{(j)} < \hat{q}_r^{(j)}} - \mathbb{1}_{X_i^{(j)} < \hat{q}_r^{(j)}} \mathbb{1}_{X_l^{(j)} < q_r^{*(j)}} \right|}_{\varepsilon_L}, \end{aligned}$$

and

$$\mathbb{E}[\varepsilon_L] \leq 2\sqrt{a_n}\mathbb{P}(\hat{q}_r^{(j)} \leq X^{(j)} < q_r^{*(j)}) + 2\sqrt{a_n}\mathbb{P}(q_r^{*(j)} \leq X^{(j)} < \hat{q}_r^{(j)}).$$

According to Lemma 1, the right hand side term converges to 0. Then, in probability,  $\lim_{n \rightarrow \infty} \varepsilon_L =$

0. Additionally,  $\lim_{n \rightarrow \infty} \frac{a_n^2}{N_{n,L}N_{n,L}^*} = \frac{1}{p_{L,r}^{(j)2}}$ , and then, in probability,

$$\lim_{n \rightarrow \infty} \sqrt{a_n}(\bar{Y}_L^* - \bar{Y}_L) = 0. \quad (\text{C.19})$$

The second term of  $a_n R_{L,r}^{(j)}$  writes

$$\begin{aligned} -a_n \times 2 \frac{N_{n,L}}{a_n} (\bar{Y}_L^* - \bar{Y}_L)(\bar{Y}_L^* - \mu) \\ = -2 \frac{N_{n,L}}{a_n} \times \sqrt{a_n}(\bar{Y}_L^* - \bar{Y}_L) \times \sqrt{a_n}(\bar{Y}_L^* - \mu), \end{aligned}$$

where in probability,  $\lim_{n \rightarrow \infty} 2 \frac{N_{n,L}}{a_n} = p_{L,r}^{(j)}$ ,  $\lim_{n \rightarrow \infty} \sqrt{a_n}(\bar{Y}_L^* - \bar{Y}_L) = 0$  according to equation C.19, and  $\sqrt{a_n}(\bar{Y}_L^* - \mu)$  converges to a normal random variable from the central limit theorem. By Slutsky theorem, in probability,  $\lim_{n \rightarrow \infty} -a_n \times 2 \frac{N_{n,L}}{a_n} (\bar{Y}_L^* - \bar{Y}_L)(\bar{Y}_L^* - \mu) = 0$ . Finally for the third term of  $a_n R_{L,r}^{(j)}$  we also use equation C.19 to conclude that in probability

$$\lim_{n \rightarrow \infty} a_n \times \frac{N_{n,L}}{a_n} (\bar{Y}_L^* - \bar{Y}_L)^2 = \lim_{n \rightarrow \infty} \frac{N_{n,L}}{a_n} [\sqrt{a_n}(\bar{Y}_L^* - \bar{Y}_L)]^2 = 0$$

Consequently,

$$\lim_{n \rightarrow \infty} a_n R_{L,r}^{(j)} = 0.$$

Symmetrically, we also have

$$\delta_R = \frac{1}{p_R} \left( \frac{1}{a_n} \sum_{i=1}^{a_n} (Y_i - \mu) \mathbb{1}_{X_i^{(j)} \geq q_r^{*(j)}} \right)^2 + R_{R,r}^{(j)},$$

with  $\lim_{n \rightarrow \infty} a_n R_{R,r}^{(j)} = 0$ , in probability.

Each component of  $\mathbf{L}_{n,\mathcal{A}_1}^{(\mathcal{C}_1)}$  writes, with  $(j, r) \in \mathcal{C}_1 \setminus (j_1, r_1)$ ,

$$\begin{aligned} L_{a_n}(\mathbb{R}^p, \hat{q}_{n,r}^{(j)}) - L_{a_n}(\mathbb{R}^p, \hat{q}_{n,r_1}^{(j_1)}) &= \frac{1}{p_{L,r}^{(j)}} \left( \frac{1}{a_n} \sum_{i=1}^{a_n} (Y_i - \mu) \mathbb{1}_{X_i^{(j)} < q_r^{*(j)}} \right)^2 \\ &+ \frac{1}{p_{R,r}^{(j)}} \left( \frac{1}{a_n} \sum_{i=1}^{a_n} (Y_i - \mu) \mathbb{1}_{X_i^{(j)} \geq q_r^{*(j)}} \right)^2 - \frac{1}{p_{L,r_1}^{(j_1)}} \left( \frac{1}{a_n} \sum_{i=1}^{a_n} (Y_i - \mu) \mathbb{1}_{X_i^{(j_1)} < q_{r_1}^{*(j_1)}} \right)^2 \\ &- \frac{1}{p_{R,r_1}^{(j_1)}} \left( \frac{1}{a_n} \sum_{i=1}^{a_n} (Y_i - \mu) \mathbb{1}_{X_i^{(j_1)} \geq q_{r_1}^{*(j_1)}} \right)^2 + R_{L,r}^{(j)} + R_{R,r}^{(j)} - R_{L,r_1}^{(j_1)} - R_{R,r_1}^{(j_1)}. \end{aligned}$$

We explicitly write  $\mathcal{C}_1 = \{(j_k, r_k)\}_{k=1, \dots, c_1}$ . Then  $\mathbf{L}_{n, \mathcal{P}_1}^{(\mathcal{C}_1)}$  can be decomposed as

$$a_n \mathbf{L}_{n, \mathcal{P}_1}^{(\mathcal{C}_1)} = h_{\mathcal{P}_1}(\mathbf{V}_n) + \mathbf{R}_{n, \mathcal{P}_1},$$

where for  $k \in \{1, \dots, c_1\}$ ,

$$\begin{aligned} V_{n, 2k-1} &= \sqrt{\frac{a_n}{p_{L, r_k}^{(j_k)}}} \frac{1}{a_n} \sum_{i=1}^{a_n} (Y_i - \mu) \mathbb{1}_{X_i^{(j_k)} < q_{r_k}^{*(j_k)}}, \\ V_{n, 2k} &= \sqrt{\frac{a_n}{p_{R, r_k}^{(j_k)}}} \frac{1}{a_n} \sum_{i=1}^{a_n} (Y_i - \mu) \mathbb{1}_{X_i^{(j_k)} \geq q_{r_k}^{*(j_k)}}. \end{aligned}$$

$h_{\mathcal{P}_1}$  is a multivariate quadratic form defined as

$$h_{\mathcal{P}_1} : \begin{pmatrix} x_1 \\ \vdots \\ x_{2c_1} \end{pmatrix} \rightarrow \begin{pmatrix} x_3^2 + x_4^2 - x_1^2 - x_2^2 \\ \vdots \\ x_{2k-1}^2 + x_{2k}^2 - x_1^2 - x_2^2 \\ \vdots \\ x_{2c_1-1}^2 + x_{2c_1}^2 - x_1^2 - x_2^2 \end{pmatrix}.$$

and  $R_{n, \mathcal{P}_1, k} = R_{L, r_k}^{(j_k)} + R_{R, r_k}^{(j_k)} - R_{L, r_1}^{(j_1)} - R_{R, r_1}^{(j_1)}$ .

From the multivariate central limit theorem,  $\mathbf{V}_n \xrightarrow[n \rightarrow \infty]{\mathcal{D}} \mathbf{V}$ , where  $\mathbf{V}$  is a gaussian vector of covariance matrix  $\text{Cov}[\mathbf{Z}]$ , and  $\mathbf{Z}$  is defined as, for  $k \in \{1, \dots, c_1\}$ ,

$$Z_{2k-1} = \frac{1}{\sqrt{p_{L, k}}} (Y - \mathbb{E}[Y]) \mathbb{1}_{X^{(j_k)} < q_{r_k}^{*(j_k)}}, Z_{2k} = \frac{1}{\sqrt{p_{R, k}}} (Y - \mathbb{E}[Y]) \mathbb{1}_{X^{(j_k)} \geq q_{r_k}^{*(j_k)}},$$

with the simplified notations  $p_{L, k} = p_{L, r_k}^{(j_k)}$  and  $p_{R, k} = p_{R, r_k}^{(j_k)}$ .

Finally, since  $\lim_{n \rightarrow \infty} \mathbf{R}_{n, \mathcal{P}_1} = \mathbf{0}$  in probability, from Slutsky's theorem and the continuous mapping theorem,  $a_n \mathbf{L}_{n, \mathcal{P}_1}^{(\mathcal{C}_1)} \xrightarrow[n \rightarrow \infty]{\mathcal{D}} h_{\mathcal{P}_1}(\mathbf{V})$ . Note that, since  $\mathbf{X}$  has a strictly positive density, each component of  $h_{\mathcal{P}_1}(\mathbf{V})$  has a strictly positive variance.  $\square$

*Proof of Lemma 5.* Consider a path  $\mathcal{P} = (j_1, r_1, \cdot)$ . Set  $\theta^{(V)} = (\theta_1^{(V)}, \cdot, \cdot) \in \Omega^{(V)}$ , a realization of the randomization of the split direction. Recalling that the best split in a random tree is the one maximizing the CART-splitting criterion, condition on  $\Theta^{(V)} = \theta^{(V)}$ ,

$$\{\mathcal{P}_1 \in T(\Theta, \mathcal{D}_n)\} = \bigcap_{(j, r) \in \theta_1^{(V)} \times \{1, \dots, q-1\} \setminus (j_1, r_1)} \{L_{a_n}(\mathbb{R}^p, \hat{q}_{n, r_1}^{(j_1)}) > L_{a_n}(\mathbb{R}^p, \hat{q}_{n, r}^{(j)})\}. \quad (\text{C.20})$$

We recall that, given  $\theta^{(V)}$ , we define the set of best theoretical cuts along the variables in  $\theta_1^{(V)}$  as

$$\mathcal{C}_1^*(\theta_1^{(V)}) = \underset{(j, r) \in \theta_1^{(V)} \times \{1, \dots, q-1\}}{\operatorname{argmax}} L^*(\mathbb{R}^p, q_r^{*(j)}).$$

Obviously if  $(j_1, r_1) \notin \theta_1^{(V)} \times \{1, \dots, q-1\}$ , the probability to select  $\mathcal{P}_1$  in the empirical and theoretical tree is null. In the sequel, we assume that  $(j_1, r_1) \in \theta_1^{(V)} \times \{1, \dots, q-1\}$  and distinguish between four cases:  $(j_1, r_1)$  is not among the best theoretical cuts  $\mathcal{C}_1^*(\theta_1^{(V)})$ , is the only element in  $\mathcal{C}_1^*(\theta_1^{(V)})$ , is one element of  $\mathcal{C}_1^*(\theta_1^{(V)})$  with a positive value of the theoretical CART-splitting criterion, or finally, is one element of  $\mathcal{C}_1^*(\theta_1^{(V)})$  that all have a null value of the theoretical CART-splitting criterion.

**Case 1** We assume that  $(j_1, r_1) \notin \mathcal{C}_1^*(\theta_1^{(V)})$ . By definition of the theoretical random forest,

$$\mathbb{P}(\mathcal{P}_1 \in T^*(\Theta) | \Theta^{(V)} = \theta^{(V)}) = 0 \quad (\text{C.21})$$

Let  $(j^*, r^*) \in \mathcal{C}_1^*(\theta_1^{(V)})$ , thus

$$\varepsilon = L^*(\mathbb{R}^p, q_{r^*}^{(j^*)}) - L^*(\mathbb{R}^p, q_{r_1}^{(j_1)}) > 0.$$

Using equation (C.20), we have:

$$\begin{aligned} \mathbb{P}(\mathcal{P}_1 \in T(\Theta, \mathcal{D}_n) | \Theta^{(V)} = \theta^{(V)}) \\ &\leq \mathbb{P}(L_{a_n}(\mathbb{R}^p, \hat{q}_{n,r_1}^{(j_1)}) > L_{a_n}(\mathbb{R}^p, \hat{q}_{n,r^*}^{(j^*)})) \\ &\leq \mathbb{P}(L_{a_n}(\mathbb{R}^p, \hat{q}_{n,r_1}^{(j_1)}) - L^*(\mathbb{R}^p, q_{r_1}^{(j_1)}) - \epsilon > L_{a_n}(\mathbb{R}^p, \hat{q}_{n,r^*}^{(j^*)}) - L^*(\mathbb{R}^p, q_{r^*}^{(j^*)})) \\ &\leq \mathbb{P}(L_{a_n}(\mathbb{R}^p, \hat{q}_{n,r_1}^{(j_1)}) - L^*(\mathbb{R}^p, q_{r_1}^{(j_1)}) \\ &\quad - (L_{a_n}(\mathbb{R}^p, \hat{q}_{n,r^*}^{(j^*)}) - L^*(\mathbb{R}^p, q_{r^*}^{(j^*)})) > \epsilon) \end{aligned}$$

Therefore, according to Lemma 3,

$$\lim_{n \rightarrow \infty} \mathbb{P}(\mathcal{P}_1 \in T(\Theta, \mathcal{D}_n) | \Theta^{(V)} = \theta^{(V)}) = 0 = \mathbb{P}(\mathcal{P}_1 \in T^*(\Theta) | \Theta^{(V)} = \theta^{(V)})$$

**Case 2** We assume that  $\mathcal{C}_1^*(\theta_1^{(V)}) = \{(j_1, r_1)\}$ . By definition of the theoretical random forest,

$$\mathbb{P}(\mathcal{P}_1 \in T^*(\Theta) | \Theta^{(V)} = \theta^{(V)}) = 1. \quad (\text{C.22})$$

Conditional on  $\Theta^{(V)} = \theta^{(V)}$ ,

$$\{\mathcal{P}_1 \in T(\Theta, \mathcal{D}_n)\}^c = \bigcup_{\substack{(j,r) \in \theta_1^{(V)} \times \{1, \dots, q-1\} \\ \setminus (j_1, r_1)}} \{L_{a_n}(\mathbb{R}^p, \hat{q}_{n,r_1}^{(j_1)}) \leq L_{a_n}(\mathbb{R}^p, \hat{q}_{n,r}^{(j)})\},$$

which leads to

$$\begin{aligned} 1 - \mathbb{P}(\mathcal{P}_1 \in T(\Theta, \mathcal{D}_n) | \Theta^{(V)} = \theta^{(V)}) \\ &\leq \sum_{(j,r) \in \theta_1^{(V)} \times \{1, \dots, q-1\} \setminus (j_1, r_1)} \mathbb{P}(L_{a_n}(\mathbb{R}^p, \hat{q}_{n,r_1}^{(j_1)}) \leq L_{a_n}(\mathbb{R}^p, \hat{q}_{n,r}^{(j)})). \end{aligned} \quad (\text{C.23})$$



From Lemma 3, for all  $j \in \theta_0^{(V)}$ ,  $r \in \{1, \dots, q-1\}$  such that  $(j, r) \neq (j_1, r_1)$ , in probability,

$$\lim_{n \rightarrow \infty} L_{a_n}(\mathbb{R}^p, \hat{q}_{n,r_1}^{(j_1)}) - L_{a_n}(\mathbb{R}^p, \hat{q}_{n,r}^{(j)}) = L^*(\mathbb{R}^p, q_{r_1}^{*(j_1)}) - L^*(\mathbb{R}^p, q_r^{*(j)}) > 0. \quad (\text{C.24})$$

Using inequality (C.23) and equation (C.24), we finally obtain,

$$\lim_{n \rightarrow \infty} \mathbb{P}(\mathcal{P}_1 \in T(\Theta, \mathcal{D}_n) | \Theta^{(V)} = \theta^{(V)}) = 1 = \mathbb{P}(\mathcal{P}_1 \in T^*(\Theta) | \Theta^{(V)} = \theta^{(V)}).$$

**Case 3** We assume that  $(j_1, r_1) \in \mathcal{C}_1^*(\theta_1^{(V)})$ ,  $|\mathcal{C}_1^*(\theta_1^{(V)})| > 1$ , and  $L^*(\mathbb{R}^p, q_{r_1}^{*(j_1)}) > 0$ . On one hand, conditional on  $\Theta^{(V)} = \theta^{(V)}$ ,

$$\{\mathcal{P}_1 \in T(\Theta, \mathcal{D}_n)\} \subset \bigcap_{(j,r) \in \mathcal{C}_1^*(\theta_1^{(V)}) \setminus (j_1, r_1)} \{L_{a_n}(\mathbb{R}^p, \hat{q}_{n,r_1}^{(j_1)}) > L_{a_n}(\mathbb{R}^p, \hat{q}_{n,r}^{(j)})\}.$$

On the other hand, conditional on  $\Theta^{(V)} = \theta^{(V)}$ ,

$$\begin{aligned} \{\mathcal{P}_1 \in T(\Theta, \mathcal{D}_n)\}^c &= \bigcup_{(j,r) \in \mathcal{C}_1^*(\theta_1^{(V)}) \setminus (j_1, r_1)} \{L_{a_n}(\mathbb{R}^p, \hat{q}_{n,r_1}^{(j_1)}) \leq L_{a_n}(\mathbb{R}^p, \hat{q}_{n,r}^{(j)})\} \\ &\quad \bigcup_{(j,r) \in \theta_1^{(V)} \times \{1, \dots, q-1\} \setminus \mathcal{C}_1^*(\theta_1^{(V)})} \{L_{a_n}(\mathbb{R}^p, \hat{q}_{n,r_1}^{(j_1)}) \leq L_{a_n}(\mathbb{R}^p, \hat{q}_{n,r}^{(j)})\}. \end{aligned}$$

Combining the two previous inclusions,

$$\begin{aligned} 0 &\leq \mathbb{P}\left(\bigcap_{(j,r) \in \mathcal{C}_1^*(\theta_1^{(V)}) \setminus (j_1, r_1)} \{L_{a_n}(\mathbb{R}^p, \hat{q}_{n,r_1}^{(j_1)}) > L_{a_n}(\mathbb{R}^p, \hat{q}_{n,r}^{(j)})\}\right) \\ &\quad - \mathbb{P}(\mathcal{P}_1 \in T(\Theta, \mathcal{D}_n) | \Theta^{(V)} = \theta^{(V)}) \\ &\leq \sum_{(j,r) \in \theta_1^{(V)} \times \{1, \dots, q-1\} \setminus \mathcal{C}_1^*(\theta_1^{(V)})} \mathbb{P}(L_{a_n}(\mathbb{R}^p, \hat{q}_{n,r_1}^{(j_1)}) \leq L_{a_n}(\mathbb{R}^p, \hat{q}_{n,r}^{(j)})). \end{aligned}$$

Using the same reasoning as in **Case 2**, we get

$$\begin{aligned} \lim_{n \rightarrow \infty} \mathbb{P}\left(\bigcap_{(j,r) \in \mathcal{C}_1^*(\theta_1^{(V)}) \setminus (j_1, r_1)} \{L_{a_n}(\mathbb{R}^p, \hat{q}_{n,r_1}^{(j_1)}) > L_{a_n}(\mathbb{R}^p, \hat{q}_{n,r}^{(j)})\}\right) \\ - \mathbb{P}(\mathcal{P}_1 \in T(\Theta, \mathcal{D}_n) | \Theta^{(V)} = \theta^{(V)}) = 0. \end{aligned}$$

We define the random vector  $\mathbf{L}_{n, \mathcal{P}_1}^{(\mathcal{C}_1^*)}$  where each component is the difference between the empirical CART-splitting criterion for the splits  $(j, r) \in \mathcal{C}_1^* \setminus (j_1, r_1)$  and  $(j_1, r_1)$ ,

$$\mathbf{L}_{n, \mathcal{P}_1}^{(\mathcal{C}_1^*)} = \left( L_{a_n}(\mathbb{R}^p, \hat{q}_{n,r}^{(j)}) - L_{a_n}(\mathbb{R}^p, \hat{q}_{n,r_1}^{(j_1)}) \right)_{(j,r) \in \mathcal{C}_1^* \setminus (j_1, r_1)},$$

then

$$\mathbb{P}\left(\bigcap_{(j,r) \in \mathcal{C}_1^*(\theta_1^{(V)}) \setminus (j_1, r_1)} \{L_{a_n}(\mathbb{R}^p, \hat{q}_{n,r_1}^{(j_1)}) > L_{a_n}(\mathbb{R}^p, \hat{q}_{n,r}^{(j)})\}\right) = \mathbb{P}(\mathbf{L}_{n, \mathcal{P}_1}^{(\mathcal{C}_1^*)} < \mathbf{0})$$

From Lemma 4 (case (a)),

$$\sqrt{a_n} \mathbf{L}_{n, \mathcal{P}_1}^{(\mathcal{C}_1^*)} \xrightarrow[n \rightarrow \infty]{\mathcal{D}} \mathcal{N}(0, \Sigma).$$

where for all  $(j, r), (j', r') \in \mathcal{C}_1^* \setminus (j_1, r_1)$ , each element of the covariance matrix  $\Sigma$  is defined by

$$\Sigma_{(j,r),(j',r')} = \text{Cov}[Z_{j,r}, Z_{j',r'}],$$

with

$$\begin{aligned} Z_{j,r} = & \left( Y - \mu_{L,r_1}^{(j_1)} \mathbb{1}_{X^{(j_1)} < q_{r_1}^{*(j_1)}} - \mu_{R,r_1}^{(j_1)} \mathbb{1}_{X^{(j_1)} \geq q_{r_1}^{*(j_1)}} \right)^2 \\ & - \left( Y - \mu_{L,r}^{(j)} \mathbb{1}_{X^{(j)} < q_r^{*(j)}} - \mu_{R,r}^{(j)} \mathbb{1}_{X^{(j)} \geq q_r^{*(j)}} \right)^2, \end{aligned}$$

$\mu_{L,r}^{(j)} = \mathbb{E}[Y|X^{(j)} < q_r^{*(j)}]$ ,  $\mu_{R,r}^{(j)} = \mathbb{E}[Y|X^{(j)} \geq q_r^{*(j)}]$ , and the variance of  $Z_{j,r}$  is strictly positive. If  $\Phi_{\theta_1^{(V)}, (j_1, r_1)}$  is the c.d.f. of the multivariate normal distribution of covariance matrix  $\Sigma$ , we can conclude

$$\begin{aligned} \lim_{n \rightarrow \infty} \mathbb{P}(\mathcal{P}_1 \in T(\Theta, \mathcal{D}_n) | \Theta^{(V)} = \theta^{(V)}) &= \lim_{n \rightarrow \infty} \mathbb{P}(\sqrt{a_n} \mathbf{L}_{n, \mathcal{P}_1}^{(\mathcal{C}_1^*)} < \mathbf{0}) \\ &= \Phi_{\theta_1^{(V)}, (j_1, r_1)}(\mathbf{0}), \end{aligned}$$

where

$$\sum_{(j,r) \in \mathcal{C}_1^*(\theta_1^{(V)})} \Phi_{\theta_1^{(V)}, (j,r)}(\mathbf{0}) = 1.$$

According to Definition 1, in the theoretical random forest, if  $\mathcal{C}_1^*(\theta_1^{(V)})$  has multiple elements,  $(j_1, r_1)$  is randomly drawn with probability

$$\mathbb{P}(\mathcal{P}_1 \in T^*(\Theta) | \Theta^{(V)} = \theta^{(V)}) = \Phi_{\theta_1^{(V)}, (j_1, r_1)}(\mathbf{0}),$$

that is

$$\begin{aligned} \lim_{n \rightarrow \infty} \mathbb{P}(\mathcal{P}_1 \in T(\Theta, \mathcal{D}_n) | \Theta^{(V)} = \theta^{(V)}) &= \mathbb{P}(\mathcal{P}_1 \in T^*(\Theta) | \Theta^{(V)} = \theta^{(V)}) \\ &= \Phi_{\theta_1^{(V)}, (j_1, r_1)}(\mathbf{0}). \end{aligned}$$

We can notice that, in the specific case where  $\mathcal{C}_1^*(\theta_1^{(V)})$  has two elements, they are both selected with equal probability  $\frac{1}{2}$ . For more than two elements, the weights are not necessary equal, it depends on the covariance matrix  $\Sigma$ .

**Case 4** We assume that all candidate splits have a null value for the theoretical CART-splitting criterion, i.e. for  $(j, r) \in \theta_1^{(V)} \times \{1, \dots, q-1\}$ ,  $L^*(\mathbb{R}^p, q_r^{*(j)}) = 0$ . Consequently  $\mathcal{C}_1^*(\theta_1^{(V)}) = \theta_1^{(V)} \times \{1, \dots, q-1\}$ . By definition

$$\mathbb{P}(\mathcal{P}_1 \in T(\Theta, \mathcal{D}_n) | \Theta^{(V)} = \theta^{(V)}) = \mathbb{P}(\mathbf{L}_{n, \mathcal{P}_1}^{(\mathcal{C}_1^*)} < \mathbf{0}).$$

According to Lemma 4 (case (b)),

$$a_n \mathbf{L}_{n, \mathcal{P}_1}^{(c_1)} \xrightarrow[n \rightarrow \infty]{\mathcal{D}} h_{\mathcal{P}_1}(\mathbf{V}),$$

where  $\mathbf{V}$  is a gaussian vector of covariance matrix  $\text{Cov}[\mathbf{Z}]$ . If  $\mathcal{C}_1^*(\theta_1^{(V)})$  is explicitly written  $\mathcal{C}_1^*(\theta_1^{(V)}) = \{(j_k, r_k)\}_{k=1, \dots, c_1}$ ,  $\mathbf{Z}$  is defined as, for  $k \in \{1, \dots, c_1\}$ ,

$$Z_{2k-1} = \frac{1}{\sqrt{p_{L,k}}}(Y - \mathbb{E}[Y])\mathbb{1}_{X^{(j_k)} < q_{r_k}^{*(j_k)}} \\ Z_{2k} = \frac{1}{\sqrt{p_{R,k}}}(Y - \mathbb{E}[Y])\mathbb{1}_{X^{(j_k)} \geq q_{r_k}^{*(j_k)}},$$

$p_{L,k} = \mathbb{P}(X^{(j_k)} < q_{r_k}^{*(j_k)})$ ,  $p_{R,k} = \mathbb{P}(X^{(j_k)} \geq q_{r_k}^{*(j_k)})$ , and  $h_{\mathcal{P}_1}$  is a multivariate quadratic form defined as

$$h_{\mathcal{P}_1} : \begin{pmatrix} x_1 \\ \vdots \\ x_{2c_1} \end{pmatrix} \rightarrow \begin{pmatrix} x_3^2 + x_4^2 - x_1^2 - x_2^2 \\ \vdots \\ x_{2k-1}^2 + x_{2k}^2 - x_1^2 - x_2^2 \\ \vdots \\ x_{2c_1-1}^2 + x_{2c_1}^2 - x_1^2 - x_2^2 \end{pmatrix}.$$

and the variance of each component of  $h_{\mathcal{P}_1}(\mathbf{V})$  is strictly positive. If  $\Phi_{\theta_1^{(V)}, (j_1, r_1)}$  is the cdf of  $h_{\mathcal{P}_1}(\mathbf{V})$ , then as in **Case 3**,

$$\begin{aligned} \lim_{n \rightarrow \infty} \mathbb{P}(\mathcal{P}_1 \in T(\Theta, \mathcal{D}_n) | \Theta^{(V)} = \theta^{(V)}) &= \Phi_{\theta_1^{(V)}, (j_1, r_1)}(\mathbf{0}) \\ &= \mathbb{P}(\mathcal{P}_1 \in T^*(\Theta) | \Theta^{(V)} = \theta^{(V)}). \end{aligned}$$

□

### C.3.2 Case 2: $\mathcal{P}_2$

*Proof of Lemma 6.* Let  $(j, r) \in \mathcal{C}_{\mathcal{P}_1}$ .

$$\begin{aligned} &\sqrt{a_n}(L_{a_n}(\hat{H}_n(\mathcal{P}_1), \hat{q}_{n,r}^{(j)}) - L_{a_n}(H^*(\mathcal{P}_1), q_r^{*(j)})) \\ &= \sqrt{a_n}[L_{a_n}(H^*(\mathcal{P}_1), \hat{q}_{n,r}^{(j)}) - L_{a_n}(H^*(\mathcal{P}_1), q_r^{*(j)})] \\ &\quad + \sqrt{a_n}[L_{a_n}(\hat{H}_n(\mathcal{P}_1), \hat{q}_{n,r}^{(j)}) - L_{a_n}(H^*(\mathcal{P}_1), \hat{q}_{n,r}^{(j)})]. \end{aligned}$$

Since  $(j, r) \in \mathcal{C}_{\mathcal{P}_1}$ ,  $\mathbb{P}(\mathbf{X} \in H^*(\mathcal{P}_1) | X^{(j)} < q_r^{*(j)}) > 0$  and  $\mathbb{P}(\mathbf{X} \in H^*(\mathcal{P}_1) | X^{(j)} \geq q_r^{*(j)}) > 0$ . Then, we can directly apply Lemma 2 to the first term of this decomposition, which shows that, in probability

$$\lim_{n \rightarrow \infty} \sqrt{a_n}(L_{a_n}(H^*(\mathcal{P}_1), \hat{q}_{n,r}^{(j)}) - L_{a_n}(H^*(\mathcal{P}_1), q_r^{*(j)})) = 0.$$

We expand the second term

$$\begin{aligned}
& \sqrt{a_n} (L_{a_n}(\hat{H}_n(\mathcal{P}_1), \hat{q}_{n,r}^{(j)}) - L_{a_n}(H^*(\mathcal{P}_1), \hat{q}_{n,r}^{(j)})) \\
&= \frac{\sqrt{a_n}}{N_n(\hat{H}_n(\mathcal{P}_1))} \sum_{i=1}^{a_n} (Y_i - \bar{Y}_{\hat{H}_n(\mathcal{P}_1)})^2 \mathbb{1}_{\mathbf{X}_i \in \hat{H}_n(\mathcal{P}_1)} \\
&- \frac{\sqrt{a_n}}{N_n(H^*(\mathcal{P}_1))} \sum_{i=1}^{a_n} (Y_i - \bar{Y}_{H^*(\mathcal{P}_1)})^2 \mathbb{1}_{\mathbf{X}_i \in H^*(\mathcal{P}_1)} \\
&- \frac{\sqrt{a_n}}{N_n(\hat{H}_n(\mathcal{P}_1))} \sum_{i=1}^{a_n} (Y_i - \bar{Y}_{\hat{H}_L} \mathbb{1}_{X_i^{(j)} < \hat{q}_{n,r}^{(j)}} - \bar{Y}_{\hat{H}_R} \mathbb{1}_{X_i^{(j)} \geq \hat{q}_{n,r}^{(j)}})^2 \mathbb{1}_{\mathbf{X}_i \in \hat{H}_n(\mathcal{P}_1)} \\
&+ \frac{\sqrt{a_n}}{N_n(H^*(\mathcal{P}_1))} \sum_{i=1}^{a_n} (Y_i - \bar{Y}_{H_L^*} \mathbb{1}_{X_i^{(j)} < \hat{q}_{n,r}^{(j)}} - \bar{Y}_{H_R^*} \mathbb{1}_{X_i^{(j)} \geq \hat{q}_{n,r}^{(j)}})^2 \mathbb{1}_{\mathbf{X}_i \in H^*(\mathcal{P}_1)}
\end{aligned}$$

with  $\hat{H}_L = \{\mathbf{x} \in \hat{H}_n(\mathcal{P}_1) : x^{(j)} < \hat{q}_{n,r}^{(j)}\}$ ,  $H_L^* = \{\mathbf{x} \in H^*(\mathcal{P}_1) : x^{(j)} < \hat{q}_{n,r}^{(j)}\}$ , and for all  $H \subseteq \mathbb{R}^p$

$$N_n(H) = \frac{1}{a_n} \sum_{i=1}^{a_n} \mathbb{1}_{\mathbf{X}_i \in H}, \quad \bar{Y}_H = \frac{1}{N_n(H)} \sum_{i=1}^{a_n} Y_i \mathbb{1}_{\mathbf{X}_i \in H}.$$

We define symmetrically  $\hat{H}_R$  and  $H_R^*$ . We obtain

$$\sqrt{a_n} (L_{a_n}(\hat{H}_n(\mathcal{P}_1), \hat{q}_{n,r}^{(j)}) - L_{a_n}(H^*(\mathcal{P}_1), \hat{q}_{n,r}^{(j)})) = \Delta_{n,1} + \Delta_{n,2} + \Delta_{n,3},$$

where

$$\Delta_{n,1} = \sqrt{a_n} (\bar{Y}_{H^*(\mathcal{P}_1)}^2 - \bar{Y}_{\hat{H}_n(\mathcal{P}_1)}^2),$$

$$\Delta_{n,2} = \sqrt{a_n} \frac{\bar{Y}_{\hat{H}_L}^2 N_n(\hat{H}_L) N_n(H^*(\mathcal{P}_1)) - \bar{Y}_{H_L^*}^2 N_n(H_L^*) N_n(\hat{H}_n(\mathcal{P}_1))}{N_n(\hat{H}_n(\mathcal{P}_1)) N_n(H^*(\mathcal{P}_1))},$$

and

$$\Delta_{n,3} = \sqrt{a_n} \frac{\bar{Y}_{\hat{H}_R}^2 N_n(\hat{H}_R) N_n(H^*(\mathcal{P}_1)) - \bar{Y}_{H_R^*}^2 N_n(H_R^*) N_n(\hat{H}_n(\mathcal{P}_1))}{N_n(\hat{H}_n(\mathcal{P}_1)) N_n(H^*(\mathcal{P}_1))}.$$

We first consider  $\Delta_{n,1}$ . Simple calculations show that

$$\begin{aligned}
\Delta_{n,1} &= \frac{\sqrt{a_n}}{N_n(H^*(\mathcal{P}_1))^2 N_n(\hat{H}_n(\mathcal{P}_1))^2} \\
&\times \sum_{i,k,l,m} Y_i Y_k [\mathbb{1}_{\mathbf{X}_i \in H^*(\mathcal{P}_1), \mathbf{X}_k \in H^*(\mathcal{P}_1), \mathbf{X}_l \in \hat{H}_n(\mathcal{P}_1), \mathbf{X}_m \in \hat{H}_n(\mathcal{P}_1)} \\
&\quad - \mathbb{1}_{\mathbf{X}_i \in \hat{H}_n(\mathcal{P}_1), \mathbf{X}_k \in \hat{H}_n(\mathcal{P}_1), \mathbf{X}_l \in H^*(\mathcal{P}_1), \mathbf{X}_m \in H^*(\mathcal{P}_1)}]
\end{aligned}$$

We consider the case  $s_1 = L$ , ( $s_1 = R$  is similar). Since  $Y_i \in \{0, 1\}$ ,

$$|\Delta_{n,1}| \leq \frac{\sqrt{a_n}}{N_n(H^*(\mathcal{P}_1))^2 N_n(\hat{H}_n(\mathcal{P}_1))^2} \times \sum_{i,k,l,m} \left| \mathbb{1}_{X_i^{(j_1)} < q_{r_1}^{*(j_1)}, X_k^{(j_1)} < q_{r_1}^{*(j_1)}, X_l^{(j_1)} < \hat{q}_{n,r_1}^{(j_1)}, X_m^{(j_1)} < \hat{q}_{n,r_1}^{(j_1)}} \right. \\ \left. - \mathbb{1}_{X_i^{(j_1)} < \hat{q}_{n,r_1}^{(j_1)}, X_k^{(j_1)} < \hat{q}_{n,r_1}^{(j_1)}, X_l^{(j_1)} < q_{r_1}^{*(j_1)}, X_m^{(j_1)} < q_{r_1}^{*(j_1)}} \right|$$

As in the proof of Lemma 2, according to Lemma 1,  $\lim_{n \rightarrow \infty} \Delta_{n,1} = 0$ , in probability. Since  $\Delta_{n,2}$  and  $\Delta_{n,3}$  are the same quantities computed on each of the two daughter nodes, we study  $\Delta_{n,2}$  only.

$$\Delta_{n,2} = \frac{\sqrt{a_n}(N_n(\hat{H}_L)N_n(H_L^*))^{-1}}{N_n(\hat{H}_n(\mathcal{P}_1))N_n(H^*(\mathcal{P}_1))} \times \sum_{i,k,l,m} Y_i Y_k \left[ \mathbb{1}_{\mathbf{X}_i \in \hat{H}_L, \mathbf{X}_k \in \hat{H}_L, \mathbf{X}_l \in H_L^*, \mathbf{X}_m \in H^*(\mathcal{P}_1)} \right. \\ \left. - \mathbb{1}_{\mathbf{X}_i \in H_L^*, \mathbf{X}_k \in H_L^*, \mathbf{X}_l \in \hat{H}_L, \mathbf{X}_m \in \hat{H}_n(\mathcal{P}_1)} \right] \\ = \frac{\sqrt{a_n}(N_n(\hat{H}_L)N_n(H_L^*))^{-1}}{N_n(\hat{H}_n(\mathcal{P}_1))N_n(H^*(\mathcal{P}_1))} \sum_{i,k,l,m} Y_i Y_k \mathbb{1}_{X_i^{(j)} < \hat{q}_{n,r}^{(j)}, X_k^{(j)} < \hat{q}_{n,r}^{(j)}, X_l^{(j)} < \hat{q}_{n,r}^{(j)}} \\ \times \left[ \mathbb{1}_{X_i^{(j_1)} < \hat{q}_{n,r_1}^{(j_1)}, X_k^{(j_1)} < \hat{q}_{n,r_1}^{(j_1)}, X_l^{(j_1)} < q_{r_1}^{*(j_1)}, X_m^{(j_1)} < q_{r_1}^{*(j_1)}} \right. \\ \left. - \mathbb{1}_{X_i^{(j_1)} < q_{r_1}^{*(j_1)}, X_k^{(j_1)} < q_{r_1}^{*(j_1)}, X_l^{(j_1)} < \hat{q}_{n,r_1}^{(j_1)}, X_m^{(j_1)} < \hat{q}_{n,r_1}^{(j_1)}} \right].$$

Therefore

$$|\Delta_{n,2}| \leq \frac{\sqrt{a_n}(N_n(\hat{H}_L)N_n(H_L^*))^{-1}}{N_n(\hat{H}_n(\mathcal{P}_1))N_n(H^*(\mathcal{P}_1))} \times \sum_{i,k,l,m} \left| \mathbb{1}_{X_i^{(j_1)} < \hat{q}_{n,r_1}^{(j_1)}, X_k^{(j_1)} < \hat{q}_{n,r_1}^{(j_1)}, X_l^{(j_1)} < q_{r_1}^{*(j_1)}, X_m^{(j_1)} < q_{r_1}^{*(j_1)}} \right. \\ \left. - \mathbb{1}_{X_i^{(j_1)} < q_{r_1}^{*(j_1)}, X_k^{(j_1)} < q_{r_1}^{*(j_1)}, X_l^{(j_1)} < \hat{q}_{n,r_1}^{(j_1)}, X_m^{(j_1)} < \hat{q}_{n,r_1}^{(j_1)}} \right|.$$

As in the proof of Lemma 2, according to Lemma 1,  $\lim_{n \rightarrow \infty} \Delta_{n,2} = 0$ , in probability, which concludes the proof, since  $\Delta_{n,3}$  can be studied in the same manner.  $\square$

*Proof of Lemma 7.* Let  $(j, r) \in \mathcal{C}_{\mathcal{P}_1}$ .

$$L_{a_n}(\hat{H}_n(\mathcal{P}_1), \hat{q}_{n,r}^{(j)}) = L_{a_n}(H^*(\mathcal{P}_1), q_r^{*(j)}) \\ + [L_{a_n}(\hat{H}_n(\mathcal{P}_1), \hat{q}_{n,r}^{(j)}) - L_{a_n}(H^*(\mathcal{P}_1), q_r^{*(j)})] \quad (\text{C.25})$$

According to Lemma 6, the second term in equation (C.25) converges to 0 in probability. From the law of large numbers, in probability,

$$\lim_{n \rightarrow \infty} L_{a_n}(H^*(\mathcal{P}_1), q_r^{*(j)}) = L^*(H^*(\mathcal{P}_1), q_r^{*(j)}),$$

which concludes the proof.  $\square$

*Proof of Lemma 8.* Similar to the case with  $\mathcal{P}_1$  (Lemma 3), where Lemma 6 is used instead of Lemma 2.  $\square$

*Proof of Lemma 9.* Consider a path  $\mathcal{P}_2 = \{(j_1, r_1, L), (j_2, r_2, \cdot)\}$ . Set  $\theta^{(V)} = (\theta_1^{(V)}, \theta_2^{(V)})$ , a realization of the randomization of the split directions at the root node and its left child node. Then,  $\theta_1^{(V)}$  and  $\theta_2^{(V)}$  denote the set of eligible variables for respectively the first and second split. We also consider  $\mathcal{C}_{\mathcal{P}_1}(\theta_2^{(V)}) \subset \mathcal{C}_{\mathcal{P}_1}$  the set of eligible second splits.

Recalling that the best split in a random tree is the one maximizing the CART-splitting criterion, conditional on  $\Theta^{(V)} = \theta^{(V)}$ ,

$$\begin{aligned} \{\mathcal{P}_2 \in T(\Theta, \mathcal{D}_n)\} = & \bigcap_{\substack{(j,r) \in \theta_1^{(V)} \times \{1, \dots, q-1\} \\ \setminus (j_1, r_1)}} \{L_{a_n}(\mathbb{R}^p, \hat{q}_{n,r_1}^{(j_1)}) > L_{a_n}(\mathbb{R}^p, \hat{q}_{n,r}^{(j)})\} \\ & \bigcap_{(j,r) \in \mathcal{C}_{\mathcal{P}_1}(\theta_2^{(V)}) \setminus (j_2, r_2)} \{L_{a_n}(\hat{H}_n(\mathcal{P}_1), \hat{q}_{n,r_2}^{(j_2)}) > L_{a_n}(\hat{H}_n(\mathcal{P}_1), \hat{q}_{n,r}^{(j)})\} \end{aligned}$$

Recall that  $\mathcal{C}_1^*(\theta_1^{(V)}) = \operatorname{argmax}_{(j,r) \in \theta_1^{(V)} \times \{1, \dots, q-1\}} L^*(\mathbb{R}^p, q_r^{*(j)})$ , and similarly

$$\mathcal{C}_2^*(\theta_2^{(V)}) = \operatorname{argmax}_{(j,r) \in \mathcal{C}_{\mathcal{P}_1}(\theta_2^{(V)})} L^*(H^*(\mathcal{P}_1), q_r^{*(j)}).$$

Obviously if  $(j_1, r_1) \notin \theta_1^{(V)} \times \{1, \dots, q-1\}$  or  $(j_2, r_2) \notin \mathcal{C}_{\mathcal{P}_1}(\theta_2^{(V)})$ , the probability to select  $\mathcal{P}_2$  in the empirical and theoretical tree is null. In the sequel, we assume that  $(j_1, r_1) \in \theta_1^{(V)} \times \{1, \dots, q-1\}$  and  $(j_2, r_2) \in \mathcal{C}_{\mathcal{P}_1}(\theta_2^{(V)})$  and distinguish between cases, depending on whether  $(j_1, r_1) \in \mathcal{C}_1^*(\theta_1^{(V)})$  or not and  $(j_2, r_2) \in \mathcal{C}_2^*(\theta_2^{(V)})$  or not, as well as the cardinality of  $\mathcal{C}_1^*(\theta_1^{(V)})$  and  $\mathcal{C}_2^*(\theta_2^{(V)})$ , and whether the maximum of the theoretical CART-splitting criterion is null or not.

**Case 1** We assume that  $(j_1, r_1) \notin \mathcal{C}_1^*(\theta_1^{(V)})$ . Hence, the theoretical decision tree satisfies

$$\mathbb{P}(\mathcal{P}_2 \in T^*(\Theta) | \Theta^{(V)} = \theta^{(V)}) = \mathbb{P}(\mathcal{P}_1 \in T^*(\Theta) | \Theta^{(V)} = \theta^{(V)}) = 0.$$

According to Lemma 5, we have

$$\begin{aligned} & \lim_{n \rightarrow \infty} \mathbb{P}(\mathcal{P}_2 \in T(\Theta, \mathcal{D}_n) | \Theta^{(V)} = \theta^{(V)}) \\ & \leq \lim_{n \rightarrow \infty} \mathbb{P}(\mathcal{P}_1 \in T(\Theta, \mathcal{D}_n) | \Theta^{(V)} = \theta^{(V)}) \\ & = 0 \\ & = \mathbb{P}(\mathcal{P}_2 \in T^*(\Theta) | \Theta^{(V)} = \theta^{(V)}). \end{aligned}$$

**Case 2** We assume that  $(j_2, r_2) \notin \mathcal{C}_2^*(\theta_2^{(V)})$ . Again, for the theoretical decision tree,

$$\mathbb{P}(\mathcal{P}_2 \in T^*(\Theta) | \Theta^{(V)} = \theta^{(V)}) = 0.$$

Letting  $(j^*, r^*) \in \mathcal{C}_2^*(\theta_2^{(V)})$ ,

$$\varepsilon = L^*(H^*(\mathcal{P}_1), q_{r^*}^{*(j^*)}) - L^*(H^*(\mathcal{P}_1), q_{r_2}^{*(j_2)}).$$

Therefore,

$$\begin{aligned} \mathbb{P}(\mathcal{P}_2 \in T(\Theta, \mathcal{D}_n) | \Theta^{(V)} = \theta^{(V)}) \\ &\leq \mathbb{P}(L_{a_n}(\hat{H}_n(\mathcal{P}_1), \hat{q}_{n,r_2}^{(j_2)}) > L_{a_n}(H^*(\mathcal{P}_1), \hat{q}_{n,r^*}^{*(j^*)})) \\ &\leq \mathbb{P}(L_{a_n}(\hat{H}_n(\mathcal{P}_1), \hat{q}_{n,r_2}^{(j_2)}) - L^*(H^*(\mathcal{P}_1), q_{r_2}^{*(j_2)}) - \epsilon \\ &\quad > L_{a_n}(\hat{H}_n(\mathcal{P}_1), \hat{q}_{n,r^*}^{*(j^*)}) - L^*(H^*(\mathcal{P}_1), q_{r^*}^{*(j^*)})) \\ &\leq \mathbb{P}(L_{a_n}(\hat{H}_n(\mathcal{P}_1), \hat{q}_{n,r_2}^{(j_2)}) - L^*(H^*(\mathcal{P}_1), q_{r_2}^{*(j_2)}) \\ &\quad - (L_{a_n}(\hat{H}_n(\mathcal{P}_1), \hat{q}_{n,r^*}^{*(j^*)}) - L^*(H^*(\mathcal{P}_1), q_{r^*}^{*(j^*)})) > \epsilon). \end{aligned}$$

Consequently, according to Lemma 7,

$$\lim_{n \rightarrow \infty} \mathbb{P}(\mathcal{P}_2 \in T(\Theta, \mathcal{D}_n) | \Theta^{(V)} = \theta^{(V)}) = 0 = \mathbb{P}(\mathcal{P}_2 \in T^*(\Theta) | \Theta^{(V)} = \theta^{(V)}).$$

**Case 3** We assume that  $(j_1, r_1) \in \mathcal{C}_1^*(\theta_1^{(V)})$  and  $\mathcal{C}_2^*(\theta_2^{(V)}) = \{(j_2, r_2)\}$ , i.e.  $(j_2, r_2)$  is the unique maximum of the theoretical CART-splitting criterion for the cell  $H^*(\mathcal{P}_1)$ . By definition of the theoretical decision tree,

$$\mathbb{P}(\mathcal{P}_2 \in T^*(\Theta) | \Theta^{(V)} = \theta^{(V)}) = \mathbb{P}(\mathcal{P}_1 \in T^*(\Theta) | \Theta^{(V)} = \theta^{(V)})$$

Conditional on  $\{\Theta^{(V)} = \theta^{(V)}\}$ ,

$$\begin{aligned} \{\mathcal{P}_2 \in T(\Theta, \mathcal{D}_n)\} &= \{\mathcal{P}_1 \in T(\Theta, \mathcal{D}_n)\} \\ &\quad \bigcap_{(j,r) \in \mathcal{C}_{\mathcal{P}_1}(\theta_2^{(V)}) \setminus (j_2, r_2)} \{L_{a_n}(\hat{H}_n(\mathcal{P}_1), \hat{q}_{n,r_2}^{(j_2)}) > L_{a_n}(\hat{H}_n(\mathcal{P}_1), \hat{q}_{n,r}^{(j)})\}. \end{aligned} \quad (\text{C.26})$$

Consequently,

$$\begin{aligned} \mathbb{P}(\mathcal{P}_2 \in T(\Theta, \mathcal{D}_n) | \Theta^{(V)} = \theta^{(V)}) \\ &\geq \mathbb{P}(\mathcal{P}_1 \in T(\Theta, \mathcal{D}_n) | \Theta^{(V)} = \theta^{(V)}) \\ &\quad - \sum_{(j,r) \in \mathcal{C}_{\mathcal{P}_1}(\theta_2^{(V)}) \setminus (j_2, r_2)} \mathbb{P}(L_{a_n}(\hat{H}_n(\mathcal{P}_1), \hat{q}_{n,r_2}^{(j_2)}) \leq L_{a_n}(\hat{H}_n(\mathcal{P}_1), \hat{q}_{n,r}^{(j)})). \end{aligned} \quad (\text{C.27})$$

For  $(j, r) \in \mathcal{C}_{\mathcal{P}_1}(\theta_2^{(V)}) \setminus (j_2, r_2)$ ,

$$L^*(H^*(\mathcal{P}_1), q_{r_2}^{*(j_2)}) - L^*(H^*(\mathcal{P}_1), q_r^{*(j)}) > 0. \quad (\text{C.28})$$

Thus, using inequalities (C.27) and (C.28), and according to Lemma 7,

$$\lim_{n \rightarrow \infty} \mathbb{P}(L_{a_n}(\hat{H}_n(\mathcal{P}_1), \hat{q}_{n,r_2}^{(j_2)}) \leq L_{a_n}(\hat{H}_n(\mathcal{P}_1), \hat{q}_{n,r}^{(j)})) = 0,$$

and thus, using (C.26) and (C.27),

$$\begin{aligned} \lim_{n \rightarrow \infty} \mathbb{P}(\mathcal{P}_2 \in T(\Theta, \mathcal{D}_n) | \Theta^{(V)} = \theta^{(V)}) \\ = \lim_{n \rightarrow \infty} \mathbb{P}(\mathcal{P}_1 \in T(\Theta, \mathcal{D}_n) | \Theta^{(V)} = \theta^{(V)}) = \mathbb{P}(\mathcal{P}_1 \in T^*(\Theta) | \Theta^{(V)} = \theta^{(V)}) \\ = \mathbb{P}(\mathcal{P}_2 \in T^*(\Theta) | \Theta^{(V)} = \theta^{(V)}), \end{aligned}$$

where the second inequality is a direct consequence of Lemma 5.

**Case 4** For the first split, we assume  $(j_1, r_1) \in \mathcal{C}_1^*(\theta_1^{(V)})$  with  $L^*(\mathbb{R}^p, q_{r_1}^{*(j_1)}) > 0$ , and for the second split,  $(j_2, r_2) \in \mathcal{C}_2^*(\theta_2^{(V)})$  with  $|\mathcal{C}_2^*(\theta_2^{(V)})| > 1$  and  $L^*(H^*(\mathcal{P}_1), q_{r_2}^{*(j_2)}) > 0$ .

On one hand, conditional on the event  $\{\Theta^{(V)} = \theta^{(V)}\}$ ,

$$\begin{aligned} \{\mathcal{P}_2 \in T(\Theta, \mathcal{D}_n)\} = & \bigcap_{\substack{(j,r) \in \theta_1^{(V)} \times \{1, \dots, q-1\} \\ \setminus (j_1, r_1)}} \{L_{a_n}(\mathbb{R}^p, \hat{q}_{n,r_1}^{(j_1)}) > L_{a_n}(\mathbb{R}^p, \hat{q}_{n,r}^{(j)})\} \\ & \bigcap_{(j,r) \in \mathcal{C}_{\mathcal{P}_1}(\theta_2^{(V)}) \setminus (j_2, r_2)} \{L_{a_n}(\hat{H}_n(\mathcal{P}_1), \hat{q}_{n,r_2}^{(j_2)}) > L_{a_n}(\hat{H}_n(\mathcal{P}_1), \hat{q}_{n,r}^{(j)})\}. \end{aligned} \quad (\text{C.29})$$

Using equation (C.29) to find a subset and a superset of  $\{\mathcal{P}_2 \in T(\Theta, \mathcal{D}_n)\}$ , we obtain

$$\begin{aligned} 0 & \geq \mathbb{P}(\mathcal{P}_2 \in T(\Theta, \mathcal{D}_n) | \Theta^{(V)} = \theta^{(V)}) \\ & - \mathbb{P}\left(\bigcap_{(j,r) \in \mathcal{C}_1^*(\theta_1^{(V)}) \setminus (j_1, r_1)} \{L_{a_n}(\mathbb{R}^p, \hat{q}_{n,r_1}^{(j_1)}) > L_{a_n}(\mathbb{R}^p, \hat{q}_{n,r}^{(j)})\} \right. \\ & \quad \left. \bigcap_{(j,r) \in \mathcal{C}_2^*(\theta_2^{(V)}) \setminus (j_2, r_2)} \{L_{a_n}(\hat{H}_n(\mathcal{P}_1), \hat{q}_{n,r_2}^{(j_2)}) > L_{a_n}(\hat{H}_n(\mathcal{P}_1), \hat{q}_{n,r}^{(j)})\} \right) \\ & \geq \sum_{(j,r) \in \theta_1^{(V)} \times \{1, \dots, q-1\} \setminus \mathcal{C}_1^*(\theta_1^{(V)})} \mathbb{P}(L_{a_n}(\mathbb{R}^p, \hat{q}_{n,r_1}^{(j_1)}) \leq L_{a_n}(\mathbb{R}^p, \hat{q}_{n,r}^{(j)})) \\ & + \sum_{(j,r) \in \theta_2^{(V)} \times \{1, \dots, q-1\} \setminus \mathcal{C}_2^*(\theta_2^{(V)})} \mathbb{P}(L_{a_n}(\hat{H}_n(\mathcal{P}_1), \hat{q}_{n,r_2}^{(j_2)}) \leq L_{a_n}(\hat{H}_n(\mathcal{P}_1), \hat{q}_{n,r}^{(j)})) \end{aligned}$$

We proved in **Case 3** that the limit of the last two terms of the previous inequality is zero, in



probability. Therefore,

$$\begin{aligned}
& \lim_{n \rightarrow \infty} \mathbb{P}(\mathcal{P}_2 \in T(\Theta, \mathcal{D}_n) | \Theta^{(V)} = \theta^{(V)}) \\
&= \lim_{n \rightarrow \infty} \mathbb{P} \left( \bigcap_{(j,r) \in \mathcal{C}_1^*(\theta_1^{(V)}) \setminus (j_1, r_1)} \{L_{a_n}(\mathbb{R}^p, \hat{q}_{n,r_1}^{(j_1)}) > L_{a_n}(\mathbb{R}^p, \hat{q}_{n,r}^{(j)})\} \right. \\
&\quad \left. \bigcap_{(j,r) \in \mathcal{C}_2^*(\theta_2^{(V)}) \setminus (j_2, r_2)} \{L_{a_n}(\hat{H}_n(\mathcal{P}_1), \hat{q}_{n,r_2}^{(j_2)}) > L_{a_n}(\hat{H}_n(\mathcal{P}_1), \hat{q}_{n,r}^{(j)})\} \right). \tag{C.30}
\end{aligned}$$

We define the random vector  $\mathbf{L}_{n, \mathcal{P}_2}^{(\mathcal{C}_1^*, \mathcal{C}_2^*)}$  (we drop  $\theta^{(V)}$  to lighten notations) where each component is the difference between the empirical CART-splitting criterion for the splits  $(j, r) \in \mathcal{C}_1^* \setminus (j_1, r_1)$  and  $(j_1, r_1)$  for the first  $|\mathcal{C}_1^*| - 1$  components, and for the splits  $(j, r) \in \mathcal{C}_2^* \setminus (j_2, r_2)$  and  $(j_2, r_2)$  for the remaining  $|\mathcal{C}_2^*| - 1$  components, i.e.,

$$\mathbf{L}_{n, \mathcal{P}_2}^{(\mathcal{C}_1^*, \mathcal{C}_2^*)} = \begin{pmatrix} [L_{a_n}(\mathbb{R}^p, \hat{q}_{n,r}^{(j)}) - L_{a_n}(\mathbb{R}^p, \hat{q}_{n,r_1}^{(j_1)})]_{(j,r) \in \mathcal{C}_1^* \setminus (j_1, r_1)} \\ [L_{a_n}(\hat{H}_n(\mathcal{P}_1), \hat{q}_{n,r}^{(j)}) - L_{a_n}(\hat{H}_n(\mathcal{P}_1), \hat{q}_{n,r_2}^{(j_2)})]_{(j,r) \in \mathcal{C}_2^* \setminus (j_2, r_2)} \end{pmatrix}.$$

Then, we can write

$$\begin{aligned}
& \mathbb{P} \left( \bigcap_{(j,r) \in \mathcal{C}_1^*(\theta_1^{(V)}) \setminus (j_1, r_1)} \{L_{a_n}(\mathbb{R}^p, \hat{q}_{n,r_1}^{(j_1)}) > L_{a_n}(\mathbb{R}^p, \hat{q}_{n,r}^{(j)})\} \right. \\
&\quad \left. \bigcap_{(j,r) \in \mathcal{C}_2^*(\theta_2^{(V)}) \setminus (j_2, r_2)} \{L_{a_n}(\hat{H}_n(\mathcal{P}_1), \hat{q}_{n,r_2}^{(j_2)}) > L_{a_n}(\hat{H}_n(\mathcal{P}_1), \hat{q}_{n,r}^{(j)})\} \right) \\
&= \mathbb{P}(\mathbf{L}_{n,2}^{(\mathcal{C}_1^*, \mathcal{C}_2^*)} < \mathbf{0}) \tag{C.31}
\end{aligned}$$

According to Lemma 8,

$$\sqrt{a_n} \mathbf{L}_{n, \mathcal{P}_2}^{(\mathcal{C}_1^*, \mathcal{C}_2^*)} \xrightarrow[n \rightarrow \infty]{\mathcal{D}} \mathcal{N}(0, \Sigma)$$

where for  $l, l' \in \{1, 2\}$ , for all  $(j, r) \in \mathcal{C}_l^* \setminus (j_l, r_l)$ ,  $(j', r') \in \mathcal{C}_{l'}^* \setminus (j_{l'}, r_{l'})$ , each element of the covariance matrix  $\Sigma$  is defined by  $\Sigma_{(j,r,l),(j',r',l')} = \text{Cov}[Z_{j,r,l}, Z_{j',r',l'}]$ , with

$$\begin{aligned}
Z_{j,r,l} &= \frac{1}{\mathbb{P}(\mathbf{X} \in H_l)} (Y - \mu_{L,r_l}^{(j_l)} \mathbf{1}_{X^{(j_l)} < q_{r_l}^{*(j_l)}} - \mu_{R,r_l}^{(j_l)} \mathbf{1}_{X^{(j_l)} \geq q_{r_l}^{*(j_l)}})^2 \mathbf{1}_{\mathbf{X} \in H_l} \\
&\quad - \frac{1}{\mathbb{P}(\mathbf{X} \in H_l)} (Y - \mu_{L,r}^{(j)} \mathbf{1}_{X^{(j)} < q_r^{*(j)}} - \mu_{R,r}^{(j)} \mathbf{1}_{X^{(j)} \geq q_r^{*(j)}})^2 \mathbf{1}_{\mathbf{X} \in H_l},
\end{aligned}$$

$\mu_{L,r}^{(j)} = \mathbb{E}[Y | X^{(j)} < q_r^{*(j)}, \mathbf{X} \in H_l]$ ,  $\mu_{R,r}^{(j)} = \mathbb{E}[Y | X^{(j)} \geq q_r^{*(j)}, \mathbf{X} \in H_l]$ , and the variance of  $Z_{j,r,l}$  is strictly positive.

Letting  $\Phi_{\mathcal{P}_1, \theta^{(V)}, (j_2, r_2)}$  be the c.d.f. of the multivariate normal distribution with covariance matrix  $\Sigma$ , and using equalities (C.30) and (C.31),

$$\lim_{n \rightarrow \infty} \mathbb{P}(\mathcal{P}_2 \in T(\Theta, \mathcal{D}_n) | \Theta^{(V)} = \theta^{(V)}) = \Phi_{\mathcal{P}_1, \theta^{(V)}, (j_2, r_2)}(\mathbf{0}).$$

We can check that

$$\sum_{(j,r) \in \mathcal{C}_2^*(\theta^{(V)})} \Phi_{\mathcal{P}_1, \theta^{(V)}, (j,r)}(\mathbf{0}) = \mathbb{P}(\mathcal{P}_1 \in T^*(\Theta) | \Theta^{(V)} = \theta^{(V)}).$$

In the theoretical random forest, the first cut  $(j_1, r_1)$  is randomly selected with probability  $\mathbb{P}(\mathcal{P}_1 \in T^*(\Theta) | \Theta^{(V)} = \theta^{(V)})$  (see the proof of Lemma 5). For the second cut, according to Definition 2, if  $\mathcal{C}_2^*(\theta_2^{(V)})$  has multiple elements,  $(j_2, r_2)$  is randomly drawn with probability

$$\frac{\Phi_{\mathcal{P}_1, \theta^{(V)}, (j_2, r_2)}(\mathbf{0})}{\mathbb{P}(\mathcal{P}_1 \in T^*(\Theta) | \Theta^{(V)} = \theta^{(V)})}$$

Since the random selection at the root node of the tree and its children nodes are independent in the theoretical algorithm,  $\mathcal{P}_2$  is selected with probability

$$\begin{aligned} \mathbb{P}(\mathcal{P}_1 \in T^*(\Theta) | \Theta^{(V)} = \theta^{(V)}) \times \frac{\Phi_{\mathcal{P}_1, \theta^{(V)}, (j_2, r_2)}(\mathbf{0})}{\mathbb{P}(\mathcal{P}_1 \in T^*(\Theta) | \Theta^{(V)} = \theta^{(V)})} \\ = \Phi_{\mathcal{P}_1, \theta^{(V)}, (j_2, r_2)}(\mathbf{0}). \end{aligned}$$

Ultimately,

$$\begin{aligned} \lim_{n \rightarrow \infty} \mathbb{P}(\mathcal{P}_2 \in T(\Theta, \mathcal{D}_n) | \Theta^{(V)} = \theta^{(V)}) &= \mathbb{P}(\mathcal{P}_2 \in T^*(\Theta) | \Theta^{(V)} = \theta^{(V)}) \\ &= \Phi_{\mathcal{P}_1, \theta^{(V)}, (j_2, r_2)}(\mathbf{0}). \end{aligned}$$

**Case 5** We assume that  $(j_1, r_1) \in \mathcal{C}_1^*(\theta_1^{(V)})$  and  $(j_2, r_2) \in \mathcal{C}_2^*(\theta_2^{(V)})$ , and that the theoretical CART-splitting criterion is null for both splits:  $L^*(\mathbb{R}^p, q_{r_1}^{*(j_1)}) = 0$  and  $L^*(H^*(\mathcal{P}_1), q_{r_2}^{*(j_2)}) = 0$ .

Consequently  $\mathcal{C}_1^*(\theta_1^{(V)}) = \theta_1^{(V)} \times \{1, \dots, q-1\}$ , and  $\mathcal{C}_2^*(\theta_2^{(V)}) = \mathcal{C}_{\mathcal{P}_1}(\theta_2^{(V)})$ . Using the same notations defined in **Case 4**, we have by definition

$$\mathbb{P}(\mathcal{P}_1 \in T(\Theta, \mathcal{D}_n) | \Theta^{(V)} = \theta^{(V)}) = \mathbb{P}(\mathbf{L}_{n, \mathcal{P}_2}^{(\mathcal{C}_1^*, \mathcal{C}_2^*)} < \mathbf{0}).$$

According to Lemma 8 (case (b)),

$$a_n \mathbf{L}_{n, \mathcal{P}_2}^{(\mathcal{C}_1^*, \mathcal{C}_2^*)} \xrightarrow[n \rightarrow \infty]{\mathcal{D}} h_{\mathcal{P}_2}(\mathbf{V}),$$

where  $\mathbf{V}$  is a gaussian vector of covariance matrix  $\text{Cov}[\mathbf{Z}]$ . If  $\mathcal{C}_1^*$  and  $\mathcal{C}_2^*$  are explicitly written  $\mathcal{C}_1^* = \{(j_k, r_k)\}_{k \in J_1}$ , and  $\mathcal{C}_2^* = \{(j_k, r_k)\}_{k \in J_2}$ , with  $J_1 = \{1, \dots, c_1 + 1\} \setminus 2$  and  $J_2 = \{2\} \cup \{c_1 + 2, \dots, c_1 + c_2\}$ ,  $\mathbf{Z}$  is defined as, for  $l \in \{1, 2\}$  and  $k \in J_l$ ,

$$\begin{aligned} Z_{2k-1} &= \frac{1}{\sqrt{p_{L,k} \mathbb{P}(\mathbf{X} \in H_l)}} (Y - \mathbb{E}[Y | \mathbf{X} \in H_l]) \mathbb{1}_{X^{(j_k)} < q_{r_k}^{*(j_k)}} \mathbb{1}_{\mathbf{X} \in H_l}, \\ Z_{2k} &= \frac{1}{\sqrt{p_{R,k} \mathbb{P}(\mathbf{X} \in H_l)}} (Y - \mathbb{E}[Y | \mathbf{X} \in H_l]) \mathbb{1}_{X^{(j_k)} \geq q_{r_k}^{*(j_k)}} \mathbb{1}_{\mathbf{X} \in H_l}, \end{aligned}$$

$p_{L,k} = \mathbb{P}(X^{(j_k)} < q_{r_k}^{*(j_k)}, \mathbf{X} \in H_l)$ ,  $p_{R,k} = \mathbb{P}(X^{(j_k)} \geq q_{r_k}^{*(j_k)}, \mathbf{X} \in H_l)$ .  $h_{\mathcal{P}_2}$  is a multivariate quadratic form defined as

$$h_{\mathcal{P}_2} : \begin{pmatrix} x_1 \\ \vdots \\ x_{2(c_1+c_2)} \end{pmatrix} \rightarrow \begin{pmatrix} x_5^2 + x_6^2 - x_1^2 - x_2^2 \\ \vdots \\ x_{2c_1+1}^2 + x_{2c_1+2}^2 - x_1^2 - x_2^2 \\ x_{2c_1+3}^2 + x_{2c_1+4}^2 - x_3^2 - x_4^2 \\ \vdots \\ x_{2(c_1+c_2)-1}^2 + x_{2(c_1+c_2)}^2 - x_3^2 - x_4^2 \end{pmatrix},$$

and the variance of each component of  $h_{\mathcal{P}_2}(\mathbf{V})$  is strictly positive.

$\Phi_{\mathcal{P}_1, \theta^{(V)}, (j_2, r_2)}$  is now defined as the cdf of  $h_{\mathcal{P}_2}(\mathbf{V})$ , and the end of the proof is identical to **Case 4**. We conclude

$$\begin{aligned} \lim_{n \rightarrow \infty} \mathbb{P}(\mathcal{P}_2 \in T(\Theta, \mathcal{D}_n) | \Theta^{(V)} = \theta^{(V)}) &= \mathbb{P}(\mathcal{P}_2 \in T^*(\Theta) | \Theta^{(V)} = \theta^{(V)}) \\ &= \Phi_{\mathcal{P}_1, \theta^{(V)}, (j_2, r_2)}(\mathbf{0}). \end{aligned}$$

**Case 6** We assume  $(j_1, r_1) \in \mathcal{C}_1^*(\theta_1^{(V)})$ ,  $(j_2, r_2) \in \mathcal{C}_2^*(\theta_2^{(V)})$  and  $|\mathcal{C}_2^*(\theta_2^{(V)})| > 1$  as in **Case 4**, but either  $L^*(\mathbb{R}^p, q_{r_1}^{*(j_1)}) = 0$  and  $L^*(H^*(\mathcal{P}_1), q_{r_2}^{*(j_2)}) > 0$ , or  $L^*(\mathbb{R}^p, q_{r_1}^{*(j_1)}) > 0$  and  $L^*(H^*(\mathcal{P}_1), q_{r_2}^{*(j_2)}) = 0$ .

The same reasoning than for **Cases 4** and **5** applies where the limit law of  $\mathbf{L}_{n, \mathcal{P}_2}^{(\mathcal{C}_1^*, \mathcal{C}_2^*)}$  has both gaussian and  $\chi$ -square components and is given by case (c) or case (d) of Lemma 8.  $\square$

## D Proof of Theorem 2

We recall Theorem 2 for the sake of clarity.

**Theorem 2.** *If  $p_0 \in [0, 1] \setminus \mathcal{U}_n$  and  $\mathcal{D}'_n = \mathcal{D}_n$ , then, conditional on  $\mathcal{D}_n$ , we have*

$$\lim_{M \rightarrow \infty} \hat{S}_{M, n, p_0} = 1 \quad \text{in probability.} \quad (\text{D.1})$$

*In addition for  $p_0 < \max \mathcal{U}_n$ ,*

$$\begin{aligned} &1 - \mathbb{E}[\hat{S}_{M, n, p_0} | \mathcal{D}_n] \\ &\stackrel{M \rightarrow \infty}{\sim} \sum_{\mathcal{P} \in \Pi} \frac{\Phi(Mp_0, M, p_n(\mathcal{P}))(1 - \Phi(Mp_0, M, p_n(\mathcal{P})))}{\frac{1}{2} \sum_{\mathcal{P}' \in \Pi} \mathbb{1}_{p_n(\mathcal{P}') > p_0} + \mathbb{1}_{p_n(\mathcal{P}') > p_0 - \rho_n(\mathcal{P}, \mathcal{P}') \frac{\sigma_n(\mathcal{P}')}{\sigma_n(\mathcal{P})} (p_0 - p_n(\mathcal{P}))}}, \end{aligned}$$

where  $\Phi(Mp_0, M, p_n(\mathcal{P}))$  is the cdf of a binomial distribution with parameter  $p_n(\mathcal{P})$ ,  $M$  trials, evaluated at  $Mp_0$ , and, for all  $\mathcal{P}, \mathcal{P}' \in \Pi$ ,

$$\sigma_n(\mathcal{P}) = \sqrt{p_n(\mathcal{P})(1 - p_n(\mathcal{P}))},$$

and

$$\rho_n(\mathcal{P}, \mathcal{P}') = \frac{\text{Cov}(\mathbb{1}_{\mathcal{P} \in T(\Theta, \mathcal{D}_n)}, \mathbb{1}_{\mathcal{P}' \in T(\Theta, \mathcal{D}_n)} | \mathcal{D}_n)}{\sigma_n(\mathcal{P})\sigma_n(\mathcal{P}')}.$$

Let  $p_0 \in [0, \max \mathcal{U}_n) \setminus \mathcal{U}_n$  and  $\mathcal{D}'_n = \mathcal{D}_n$ . Before proving Theorem 2, we need the following two lemmas.

**Lemma 10.** *Let  $F$  be the hypergeometric function. Then, for  $(a, c) \in \mathbb{Z}^2$  and  $\mathcal{P} \in \Pi$  such that  $p_n(\mathcal{P}) > p_0$ , we have*

$$\lim_{M \rightarrow \infty} \frac{F(M + a, 1, M(1 - p_0) + c, 1 - p_n(\mathcal{P}))}{F(M + 1, 1, M(1 - p_0) + 1, 1 - p_n(\mathcal{P}))} = 1.$$

**Lemma 11.** *Let  $\mathcal{P}' \in \Pi$ . For all  $\mathcal{P} \in \Pi$  such that  $p_n(\mathcal{P}) > p_0$ , we have*

$$\lim_{M \rightarrow \infty} \mathbb{P}(\hat{p}_{M,n}(\mathcal{P}') > p_0 | \hat{p}_{M,n}(\mathcal{P}) > p_0, \mathcal{D}_n) = \mathbb{1}_{p_n(\mathcal{P}') > p_0}$$

and

$$\lim_{M \rightarrow \infty} \mathbb{P}(\hat{p}_{M,n}(\mathcal{P}') > p_0 | \hat{p}_{M,n}(\mathcal{P}) \leq p_0, \mathcal{D}_n) = \mathbb{1}_{p_n(\mathcal{P}') > p_0 - \rho_n(\mathcal{P}, \mathcal{P}') \frac{\sigma_n(\mathcal{P}')}{\sigma_n(\mathcal{P})} \times (p_0 - p_n(\mathcal{P}))}.$$

Symmetrically, for all  $\mathcal{P} \in \Pi$  such that  $p_n(\mathcal{P}) \leq p_0$ , we have

$$\begin{aligned} \lim_{M \rightarrow \infty} \mathbb{P}(\hat{p}_{M,n}(\mathcal{P}') > p_0 | \hat{p}_{M,n}(\mathcal{P}) \leq p_0, \mathcal{D}_n) &= \mathbb{1}_{p_n(\mathcal{P}') > p_0}, \\ \lim_{M \rightarrow \infty} \mathbb{P}(\hat{p}_{M,n}(\mathcal{P}') > p_0 | \hat{p}_{M,n}(\mathcal{P}) > p_0, \mathcal{D}_n) &= \mathbb{1}_{p_n(\mathcal{P}') > p_0 - \rho_n(\mathcal{P}, \mathcal{P}') \frac{\sigma_n(\mathcal{P}')}{\sigma_n(\mathcal{P})} \times (p_0 - p_n(\mathcal{P}))}. \end{aligned}$$

We are now in a position to prove Theorem 2.

*Proof of Theorem 2.* The first statement, identity (D.1), is proved similarly to Corollary 2, using the law of large numbers instead of Theorem 1. For the second statement, we first recall that, by definition,

$$\begin{aligned} \hat{S}_{M,n,p_0} &= \frac{2 \sum_{\mathcal{P} \in \Pi} \mathbb{1}_{\hat{p}_{M,n}(\mathcal{P}) > p_0} \mathbb{1}_{\hat{p}'_{M,n}(\mathcal{P}) > p_0}}{\sum_{\mathcal{P} \in \Pi} \mathbb{1}_{\hat{p}_{M,n}(\mathcal{P}) > p_0} + \mathbb{1}_{\hat{p}'_{M,n}(\mathcal{P}) > p_0}} \\ &= 1 - \frac{\sum_{\mathcal{P} \in \Pi} \mathbb{1}_{\hat{p}_{M,n}(\mathcal{P}) > p_0} \mathbb{1}_{\hat{p}'_{M,n}(\mathcal{P}) \leq p_0} + \mathbb{1}_{\hat{p}_{M,n}(\mathcal{P}) \leq p_0} \mathbb{1}_{\hat{p}'_{M,n}(\mathcal{P}) > p_0}}{\sum_{\mathcal{P} \in \Pi} \mathbb{1}_{\hat{p}_{M,n}(\mathcal{P}) > p_0} + \mathbb{1}_{\hat{p}'_{M,n}(\mathcal{P}) > p_0}}. \end{aligned}$$

Taking the expectation conditional on  $\mathcal{D}_n$  gives

$$\begin{aligned} \mathbb{E}[\hat{S}_{M,n,p_0} | \mathcal{D}_n] &= 1 - 2 \mathbb{E} \left[ \frac{\sum_{\mathcal{P} \in \Pi} \mathbb{1}_{\hat{p}_{M,n}(\mathcal{P}) > p_0} \mathbb{1}_{\hat{p}'_{M,n}(\mathcal{P}) \leq p_0}}{\sum_{\mathcal{P} \in \Pi} \mathbb{1}_{\hat{p}_{M,n}(\mathcal{P}) > p_0} + \mathbb{1}_{\hat{p}'_{M,n}(\mathcal{P}) > p_0}} \middle| \mathcal{D}_n \right] \\ &= 1 - 2 \mathbb{E} \left[ \frac{U_M}{V_M + V'_M} \middle| \mathcal{D}_n \right], \end{aligned}$$

where  $U_M = \sum_{\mathcal{P} \in \Pi} \mathbf{1}_{\hat{p}_{M,n}(\mathcal{P}) > p_0 \cap \hat{p}'_{M,n}(\mathcal{P}) \leq p_0}$ ,  $V_M = \sum_{\mathcal{P} \in \Pi} \mathbf{1}_{\hat{p}_{M,n}(\mathcal{P}) > p_0}$ , and  $V'_M = \sum_{\mathcal{P} \in \Pi} \mathbf{1}_{\hat{p}'_{M,n}(\mathcal{P}) > p_0}$ . Note that

$$\begin{aligned}\mathbb{E}[V_M | \mathcal{D}_n] &= \sum_{\mathcal{P} \in \Pi} \mathbb{P}(\hat{p}_{M,n}(\mathcal{P}) > p_0 | \mathcal{D}_n) \xrightarrow{M \rightarrow \infty} \sum_{\mathcal{P} \in \Pi} \mathbf{1}_{p_n(\mathcal{P}) > p_0}, \\ \mathbb{E}[U_M | \mathcal{D}_n] &= \sum_{\mathcal{P} \in \Pi} \mathbb{P}(\hat{p}_{M,n}(\mathcal{P}) > p_0 | \mathcal{D}_n) \mathbb{P}(\hat{p}'_{M,n}(\mathcal{P}) \leq p_0 | \mathcal{D}_n) \xrightarrow{M \rightarrow \infty} 0.\end{aligned}$$

Also,

$$\begin{aligned}\mathbb{E}\left[\frac{U_M}{V_M + V'_M} | \mathcal{D}_n\right] &= \sum_{m, m'} \frac{1}{m + m'} \mathbb{E}[U_M | V_M = m, V'_M = m', \mathcal{D}_n] \\ &\quad \times \mathbb{P}(V_M = m | \mathcal{D}_n) \mathbb{P}(V'_M = m' | \mathcal{D}_n) \\ &= \sum_{m, m'} \frac{1}{m + m'} \mathbb{E}\left[\sum_{\mathcal{P} \in \Pi} \mathbf{1}_{\hat{p}_{M,n}(\mathcal{P}) > p_0 \cap \hat{p}'_{M,n}(\mathcal{P}) \leq p_0} | V_M = m, V'_M = m', \mathcal{D}_n\right] \\ &\quad \times \mathbb{P}(V_M = m | \mathcal{D}_n) \mathbb{P}(V'_M = m' | \mathcal{D}_n) \\ &= \sum_{m, m'} \frac{1}{m + m'} \sum_{\mathcal{P} \in \Pi} \mathbb{P}(\hat{p}_{M,n}(\mathcal{P}) > p_0 | V_M = m, \mathcal{D}_n) \\ &\quad \times \mathbb{P}(\hat{p}'_{M,n}(\mathcal{P}) \leq p_0 | V'_M = m', \mathcal{D}_n) \mathbb{P}(V_M = m | \mathcal{D}_n) \mathbb{P}(V'_M = m' | \mathcal{D}_n) \\ &= \sum_{m, m'} \frac{1}{m + m'} \sum_{\mathcal{P} \in \Pi} \mathbb{P}(\hat{p}_{M,n}(\mathcal{P}) > p_0, V_M = m | \mathcal{D}_n) \\ &\quad \times \mathbb{P}(\hat{p}'_{M,n}(\mathcal{P}) \leq p_0, V'_M = m' | \mathcal{D}_n) \\ &= \sum_{\mathcal{P} \in \Pi} \mathbb{P}(\hat{p}_{M,n}(\mathcal{P}) > p_0 | \mathcal{D}_n) \mathbb{P}(\hat{p}'_{M,n}(\mathcal{P}) \leq p_0 | \mathcal{D}_n) \\ &\quad \times \left[ \sum_{m, m'} \frac{1}{m + m'} \mathbb{P}(V_M = m | \hat{p}_{M,n}(\mathcal{P}) > p_0, \mathcal{D}_n) \right. \\ &\quad \left. \times \mathbb{P}(V'_M = m' | \hat{p}'_{M,n}(\mathcal{P}) \leq p_0, \mathcal{D}_n) \right].\end{aligned}$$

For all  $\mathcal{P} \in \Pi$ ,

$$\begin{aligned}\mathbb{P}(\hat{p}_{M,n}(\mathcal{P}) > p_0 | \mathcal{D}_n) \mathbb{P}(\hat{p}'_{M,n}(\mathcal{P}) \leq p_0 | \mathcal{D}_n) \\ = \Phi(Mp_0, M, p_n(\mathcal{P}))(1 - \Phi(Mp_0, M, p_n(\mathcal{P}))),\end{aligned}$$

where  $\Phi$  is the cdf of the binomial distribution. As a direct consequence of Lemma 11,

$$\begin{aligned}\lim_{M \rightarrow \infty} \sum_{m, m'} \frac{1}{m + m'} \mathbb{P}(V_M = m | \hat{p}_{M,n}(\mathcal{P}) > p_0, \mathcal{D}_n) \\ \times \mathbb{P}(V'_M = m' | \hat{p}'_{M,n}(\mathcal{P}) \leq p_0, \mathcal{D}_n) \\ = \frac{1}{\sum_{\mathcal{P}' \in \Pi} \mathbf{1}_{p_n(\mathcal{P}') > p_0} + \mathbf{1}_{p_n(\mathcal{P}') + \rho_n(\mathcal{P}, \mathcal{P}') \frac{\sigma_n(\mathcal{P}')}{\sigma_n(\mathcal{P})} (p_0 - p_n(\mathcal{P})) > p_0}},\end{aligned}$$

which yields

$$1 - \mathbb{E}[\hat{S}_{M,n,p_0} | \mathcal{D}_n] \\ \stackrel{M \rightarrow \infty}{\sim} \sum_{\mathcal{P} \in \Pi} \frac{2\Phi(Mp_0, M, p_n(\mathcal{P}))(1 - \Phi(Mp_0, M, p_n(\mathcal{P})))}{\sum_{\mathcal{P}' \in \Pi} \mathbb{1}_{\hat{p}_n(\mathcal{P}') > p_0} + \mathbb{1}_{p_n(\mathcal{P}') + \rho_n(\mathcal{P}, \mathcal{P}') \frac{\sigma_n(\mathcal{P}')}{\sigma_n(\mathcal{P})} (p_0 - p_n(\mathcal{P})) > p_0}}.$$

This is the desired result.  $\square$

## D.1 Proof of intermediate lemmas

*Proof of lemma 10.* Cvitković et al. (2017) provides an asymptotic expansion of the hypergeometric function  $F$  in the case where the first and third parameters goes to infinity with a constant ratio. For  $a, c, z, \varepsilon \in \mathbb{R}$ ,  $b \notin \mathbb{Z} \setminus \mathbb{N}$ , such that  $\varepsilon > 1$ , and  $z\varepsilon < 1$ , Cvitković et al. (2017) gives in the section 2.2.2 (end of page 10)

$$F(a + \varepsilon\lambda, b, c + \lambda, z) \underset{|\lambda| \rightarrow \infty}{\sim} \frac{1}{(1 - \varepsilon z)^b}. \quad (\text{D.2})$$

We can then derive the limit of the following ratio

$$\lim_{|\lambda| \rightarrow \infty} \frac{F(a + \varepsilon\lambda, b, c + \lambda, z)}{F(1 + \varepsilon\lambda, b, 1 + \lambda, z)} = 1 \quad (\text{D.3})$$

We use D.3 in the specific case where  $b = 1$ ,  $a, c \in \mathbb{Z}$ ,  $\varepsilon = \frac{1}{1-p_0} > 1$ ,  $z = 1 - p_n(\mathcal{P})$  for  $\mathcal{P} \in \Pi$  such that  $p_n(\mathcal{P}) > p_0$  (and then  $z\varepsilon < 1$ ), and  $\lambda = M(1 - p_0)$ , it follows that

$$\lim_{M \rightarrow \infty} \frac{F(M + a, 1, M(1 - p_0) + c, 1 - p_n(\mathcal{P}))}{F(M + 1, 1, M(1 - p_0) + 1, 1 - p_n(\mathcal{P}))} = 1 \quad (\text{D.4})$$

$\square$

*Proof of lemma 11.* Fix  $\mathcal{D}_n$ . Let  $\mathcal{P}', \mathcal{P} \in \Pi$ . In what follows, when there is no ambiguity, we will replace  $T(\Theta, \mathcal{D}_n)$  by  $T_n(\Theta)$  to lighten notations.

**Case 1:**  $p_n(\mathcal{P}) > p_0$

$$\begin{aligned} & \mathbb{E}[\hat{p}_{M,n}(\mathcal{P}') | \hat{p}_{M,n}(\mathcal{P}) \leq p_0, \mathcal{D}_n] \\ &= \mathbb{E}\left[\frac{1}{M} \sum_{l=1}^M \mathbb{1}_{\mathcal{P}' \in T_n(\Theta_l)} \mid \hat{p}_{M,n}(\mathcal{P}) \leq p_0, \mathcal{D}_n\right] \\ &= \mathbb{P}(\mathcal{P}' \in T_n(\Theta_1) | \mathcal{P} \in T_n(\Theta_1), \hat{p}_{M,n}(\mathcal{P}) \leq p_0, \mathcal{D}_n) \\ & \quad \times \mathbb{P}(\mathcal{P} \in T_n(\Theta_1) | \hat{p}_{M,n}(\mathcal{P}) \leq p_0, \mathcal{D}_n) \\ & \quad + \mathbb{P}(\mathcal{P}' \in T_n(\Theta_1) | \mathcal{P} \notin T_n(\Theta_1), \hat{p}_{M,n}(\mathcal{P}) \leq p_0, \mathcal{D}_n) \\ & \quad \times (1 - \mathbb{P}(\mathcal{P} \in T_n(\Theta_1) | \hat{p}_{M,n}(\mathcal{P}) \leq p_0, \mathcal{D}_n)) \\ &= \mathbb{P}(\mathcal{P}' \in T_n(\Theta_1) | \mathcal{P} \in T_n(\Theta_1), \mathcal{D}_n) \mathbb{P}(\mathcal{P} \in T_n(\Theta_1) | \hat{p}_{M,n}(\mathcal{P}) \leq p_0, \mathcal{D}_n) \\ & \quad + \mathbb{P}(\mathcal{P}' \in T_n(\Theta_1) | \mathcal{P} \notin T_n(\Theta_1), \mathcal{D}_n) \\ & \quad \times (1 - \mathbb{P}(\mathcal{P} \in T_n(\Theta_1) | \hat{p}_{M,n}(\mathcal{P}) \leq p_0, \mathcal{D}_n)). \end{aligned} \quad (\text{D.5})$$

since

$$\begin{aligned} & \mathbb{P}(\mathcal{P}' \in T_n(\Theta_1) | \mathcal{P} \in T_n(\Theta_1), \hat{p}_{M,n}(\mathcal{P}) \leq p_0, \mathcal{D}_n) \\ &= \mathbb{P}(\mathcal{P}' \in T_n(\Theta_1) | \mathcal{P} \in T_n(\Theta_1), \mathcal{D}_n). \end{aligned} \quad (\text{D.6})$$

because, conditional on  $\mathcal{D}_n$ , the events  $\mathcal{P}' \in T_n(\Theta_1), \dots, \mathcal{P}' \in T_n(\Theta_M)$  are independent. We can rewrite,

$$\begin{aligned} & \mathbb{P}(\mathcal{P}' \in T_n(\Theta_1) | \mathcal{P} \notin T_n(\Theta_1), \mathcal{D}_n) \\ &= \frac{\mathbb{P}(\mathcal{P}' \in T_n(\Theta_1), \mathcal{P} \notin T_n(\Theta_1) | \mathcal{D}_n)}{1 - p_n(\mathcal{P})} \\ &= \frac{(1 - \mathbb{P}(\mathcal{P} \in T_n(\Theta_1) | \mathcal{P}' \in T_n(\Theta_1), \mathcal{D}_n)) p_n(\mathcal{P}')}{1 - p_n(\mathcal{P})} \\ &= \frac{p_n(\mathcal{P}')}{1 - p_n(\mathcal{P})} - \frac{p_n(\mathcal{P})}{1 - p_n(\mathcal{P})} \mathbb{P}(\mathcal{P}' \in T_n(\Theta_1) | \mathcal{P} \in T_n(\Theta_1), \mathcal{D}_n), \end{aligned} \quad (\text{D.7})$$

yielding, using equation (D.5),

$$\begin{aligned} & \mathbb{E}[\hat{p}_{M,n}(\mathcal{P}') | \hat{p}_{M,n}(\mathcal{P}) \leq p_0, \mathcal{D}_n] \\ &= \mathbb{P}(\mathcal{P}' \in T_n(\Theta_1) | \mathcal{P} \in T_n(\Theta_1), \mathcal{D}_n) \left( \frac{\mathbb{P}(\mathcal{P} \in T_n(\Theta_1) | \hat{p}_{M,n}(\mathcal{P}) \leq p_0, \mathcal{D}_n)}{1 - p_n(\mathcal{P})} \right. \\ & \quad \left. - \frac{p_n(\mathcal{P})}{1 - p_n(\mathcal{P})} \right) + \frac{p_n(\mathcal{P}')}{1 - p_n(\mathcal{P})} (1 - \mathbb{P}(\mathcal{P} \in T_n(\Theta_1) | \hat{p}_{M,n}(\mathcal{P}) \leq p_0, \mathcal{D}_n)). \end{aligned} \quad (\text{D.8})$$

Besides, by definition of the correlation

$$\rho_n(\mathcal{P}, \mathcal{P}') = \frac{\text{Cov}(\mathbb{1}_{\mathcal{P} \in T_n(\Theta)}, \mathbb{1}_{\mathcal{P}' \in T_n(\Theta)} | \mathcal{D}_n)}{\sigma_n(\mathcal{P}) \sigma_n(\mathcal{P}')},$$

simple calculations show that

$$\begin{aligned} & \mathbb{P}(\mathcal{P}' \in T_n(\Theta_1) | \mathcal{P} \in T_n(\Theta_1), \mathcal{D}_n) \\ &= p_n(\mathcal{P}') + \rho_n(\mathcal{P}, \mathcal{P}') \sqrt{\frac{p_n(\mathcal{P}')}{p_n(\mathcal{P})} (1 - p_n(\mathcal{P})) (1 - p_n(\mathcal{P}'))}, \end{aligned} \quad (\text{D.9})$$

which, together with equation (D.8) leads to,

$$\begin{aligned} & \mathbb{E}[\hat{p}_{M,n}(\mathcal{P}') | \hat{p}_{M,n}(\mathcal{P}) \leq p_0, \mathcal{D}_n] \\ &= p_n(\mathcal{P}') + \rho_n(\mathcal{P}, \mathcal{P}') \frac{\sigma_n(\mathcal{P}')}{\sigma_n(\mathcal{P})} (\mathbb{P}(\mathcal{P} \in T_n(\Theta_1) | \hat{p}_{M,n}(\mathcal{P}) \leq p_0, \mathcal{D}_n) \\ & \quad - p_n(\mathcal{P})). \end{aligned} \quad (\text{D.10})$$

Regarding the probability in the right-hand side of equation (D.10), we have

$$\begin{aligned}
& \mathbb{P}(\mathcal{S} \in T_n(\Theta_1) | \hat{p}_{M,n}(\mathcal{S}) \leq p_0, \mathcal{D}_n) \\
&= p_n(\mathcal{S}) \frac{\mathbb{P}(\hat{p}_{M,n}(\mathcal{S}) \leq p_0 | \mathcal{S} \in T_n(\Theta_1), \mathcal{D}_n)}{\mathbb{P}(\hat{p}_{M,n}(\mathcal{S}) \leq p_0 | \mathcal{D}_n)} \\
&= p_n(\mathcal{S}) \frac{\mathbb{P}((M-1)\hat{p}_{M-1,n}(\mathcal{S}) \leq Mp_0 - 1 | \mathcal{D}_n)}{\mathbb{P}(M\hat{p}_{M,n}(\mathcal{S}) \leq Mp_0 | \mathcal{D}_n)} \\
&= p_n(\mathcal{S}) \frac{\Phi(Mp_0 - 1, M-1, p_n(\mathcal{S}))}{\Phi(Mp_0, M, p_n(\mathcal{S}))}.
\end{aligned}$$

Using standard formulas,  $\Phi$  can be expressed with the incomplete beta function,

$$\Phi(k, M, p) = I_{1-p}(M-k, k+1) = \frac{B_{1-p}(M-k, k+1)}{B(M-k, k+1)},$$

and the regularized beta function is related to the hypergeometric function  $F$ , for  $a > 0$ ,  $b > 0$ , and  $p \in [0, 1]$  (Olver et al., 2010),

$$B_{1-p}(a, b) = \frac{(1-p)^a p^b}{a} F(a+b, 1, a+1, 1-p).$$

Then, we can express the cdf of the binomial distribution using the hypergeometric function, and it follows

$$\begin{aligned}
& \mathbb{P}(\mathcal{S} \in T_n(\Theta_1) | \hat{p}_{M,n}(\mathcal{S}) \leq p_0, \mathcal{D}_n) \\
&= p_0 \frac{F(M, 1, M(1-p_0) + 1, 1 - \hat{p}_n(\mathcal{S}))}{F(M+1, 1, M(1-p_0) + 1, 1 - \hat{p}_n(\mathcal{S}))}
\end{aligned} \tag{D.11}$$

According to Lemma 10,

$$\lim_{M \rightarrow \infty} \frac{F(M, 1, M(1-p_0) + 1, 1 - p_n(\mathcal{S}))}{F(M+1, 1, M(1-p_0) + 1, 1 - p_n(\mathcal{S}))} = 1. \tag{D.12}$$

Consequently,

$$\lim_{M \rightarrow \infty} \mathbb{P}(\mathcal{S} \in T(\Theta_1, \mathcal{D}_n) | \hat{p}_{M,n}(\mathcal{S}) \leq p_0, \mathcal{D}_n) = p_0, \tag{D.13}$$

and using this limiting result with equation (D.10) yields,

$$\begin{aligned}
\lim_{M \rightarrow \infty} \mathbb{E}[\hat{p}_{M,n}(\mathcal{S}') | \hat{p}_{M,n}(\mathcal{S}) \leq p_0, \mathcal{D}_n] &= p_n(\mathcal{S}') + \rho_n(\mathcal{S}, \mathcal{S}') \frac{\sigma_n(\mathcal{S}')}{\sigma_n(\mathcal{S})} \\
&\quad \times (p_0 - p_n(\mathcal{S})).
\end{aligned} \tag{D.14}$$

Regarding the conditional variance,

$$\begin{aligned}
& \mathbb{V}[\hat{p}_{M,n}(\mathcal{S}') | \hat{p}_{M,n}(\mathcal{S}) \leq p_0, \mathcal{D}_n] \\
&= \mathbb{V}\left[\frac{1}{M} \sum_{l=1}^M \mathbf{1}_{\mathcal{S}' \in T_n(\Theta_l)} | \hat{p}_{M,n}(\mathcal{S}) \leq p_0, \mathcal{D}_n\right]
\end{aligned}$$



$$\begin{aligned}
& \mathbb{V}[\hat{p}_{M,n}(\mathcal{P}') | \hat{p}_{M,n}(\mathcal{P}) \leq p_0, \mathcal{D}_n] \\
&= \frac{1}{M} \mathbb{V}[\mathbf{1}_{\mathcal{P}' \in T(\Theta_1, \mathcal{D}_n)} | \hat{p}_{M,n}(\mathcal{P}) \leq p_0, \mathcal{D}_n] \\
&\quad + (1 - \frac{1}{M}) \text{Cov}(\mathbf{1}_{\mathcal{P}' \in T_n(\Theta_1)}, \mathbf{1}_{\mathcal{P}' \in T_n(\Theta_2)} | \hat{p}_{M,n}(\mathcal{P}) \leq p_0, \mathcal{D}_n) \\
&\leq \frac{1}{M} + C_M
\end{aligned}$$

where

$$\begin{aligned}
C_M &= \text{Cov}(\mathbf{1}_{\mathcal{P}' \in T_n(\Theta_1)}, \mathbf{1}_{\mathcal{P}' \in T_n(\Theta_2)} | \hat{p}_{M,n}(\mathcal{P}) \leq p_0, \mathcal{D}_n) \\
&= \mathbb{P}(\mathcal{P}' \in T_n(\Theta_1), \mathcal{P}' \in T_n(\Theta_2) | \hat{p}_{M,n}(\mathcal{P}) \leq p_0, \mathcal{D}_n) \\
&\quad - \mathbb{P}(\mathcal{P}' \in T_n(\Theta_1) | \hat{p}_{M,n}(\mathcal{P}) \leq p_0, \mathcal{D}_n) \\
&\quad \times \mathbb{P}(\mathcal{P}' \in T_n(\Theta_2) | \hat{p}_{M,n}(\mathcal{P}) \leq p_0, \mathcal{D}_n)
\end{aligned}$$

Then, we follow the same reasoning that leads to equation (D.13). We can fully expand  $C_M$  using Bayes formula, depending whether  $\mathcal{P} \in T_n(\Theta_1)$  or  $\mathcal{P} \in T_n(\Theta_2)$ . Note that, since all the trees are independent conditional on  $\mathcal{D}_n$ , we can reduce the conditioning event  $\{\mathcal{P} \in T_n(\Theta_1), \mathcal{P} \in T_n(\Theta_2), \hat{p}_{M,n}(\mathcal{P}) \leq p_0, \mathcal{D}_n\}$  to  $\{\mathcal{P} \in T_n(\Theta_1), \mathcal{P} \in T_n(\Theta_2), \mathcal{D}_n\}$ , then

$$\begin{aligned}
C_M &= \mathbb{P}(\mathcal{P}' \in T_n(\Theta_1), \mathcal{P}' \in T_n(\Theta_2) | \mathcal{P} \in T_n(\Theta_1), \mathcal{P} \in T_n(\Theta_2), \mathcal{D}_n) \\
&\quad \times \mathbb{P}(\mathcal{P} \in T_n(\Theta_1), \mathcal{P} \in T_n(\Theta_2) | \hat{p}_{M,n}(\mathcal{P}) \leq p_0, \mathcal{D}_n) \\
&\quad - (\mathbb{P}(\mathcal{P}' \in T_n(\Theta_1) | \mathcal{P} \in T_n(\Theta_1), \mathcal{D}_n) \\
&\quad \times \mathbb{P}(\mathcal{P} \in T_n(\Theta_1) | \hat{p}_{M,n}(\mathcal{P}) \leq p_0, \mathcal{D}_n))^2 \\
&+ 2[\mathbb{P}(\mathcal{P}' \in T_n(\Theta_1), \mathcal{P}' \in T_n(\Theta_2) | \mathcal{P} \in T_n(\Theta_1), \mathcal{P} \notin T_n(\Theta_2), \mathcal{D}_n) \\
&\quad \times \mathbb{P}(\mathcal{P} \in T_n(\Theta_1), \mathcal{P} \notin T_n(\Theta_2) | \hat{p}_{M,n}(\mathcal{P}) \leq p_0, \mathcal{D}_n) \\
&\quad - \mathbb{P}(\mathcal{P}' \in T_n(\Theta_1) | \mathcal{P} \in T_n(\Theta_1), \mathcal{D}_n) \\
&\quad \times \mathbb{P}(\mathcal{P} \in T_n(\Theta_1) | \hat{p}_{M,n}(\mathcal{P}) \leq p_0, \mathcal{D}_n) \\
&\quad \times \mathbb{P}(\mathcal{P}' \in T_n(\Theta_1) | \mathcal{P} \notin T_n(\Theta_1), \mathcal{D}_n) \\
&\quad \times \mathbb{P}(\mathcal{P} \notin T_n(\Theta_1) | \hat{p}_{M,n}(\mathcal{P}) \leq p_0, \mathcal{D}_n)] \\
&+ \mathbb{P}(\mathcal{P}' \in T_n(\Theta_1), \mathcal{P}' \in T_n(\Theta_2) | \mathcal{P} \notin T_n(\Theta_1), \mathcal{P} \notin T_n(\Theta_2), \mathcal{D}_n) \\
&\quad \times \mathbb{P}(\mathcal{P} \notin T_n(\Theta_1), \mathcal{P} \notin T_n(\Theta_2) | \hat{p}_{M,n}(\mathcal{P}) \leq p_0, \mathcal{D}_n) \\
&\quad - (\mathbb{P}(\mathcal{P}' \in T_n(\Theta_1) | \mathcal{P} \notin T_n(\Theta_1), \mathcal{D}_n) \\
&\quad \times \mathbb{P}(\mathcal{P} \notin T_n(\Theta_1) | \hat{p}_{M,n}(\mathcal{P}) \leq p_0, \mathcal{D}_n))^2
\end{aligned}$$

Conditional on  $\mathcal{D}_n$ ,  $T_n(\Theta_1)$  and  $T_n(\Theta_2)$  are independent, then

$$\begin{aligned}
& \mathbb{P}(\mathcal{P}' \in T_n(\Theta_1), \mathcal{P}' \in T_n(\Theta_2) | \mathcal{P} \in T_n(\Theta_1), \mathcal{P} \in T_n(\Theta_2), \mathcal{D}_n) \\
&= \frac{\mathbb{P}(\mathcal{P}' \in T_n(\Theta_1), \mathcal{P}' \in T_n(\Theta_2), \mathcal{P} \in T_n(\Theta_1), \mathcal{P} \in T_n(\Theta_2) | \mathcal{D}_n)}{\mathbb{P}(\mathcal{P} \in T_n(\Theta_1), \mathcal{P} \in T_n(\Theta_2) | \mathcal{D}_n)} \\
&= \frac{\mathbb{P}(\mathcal{P}' \in T_n(\Theta_1), \mathcal{P} \in T_n(\Theta_1) | \mathcal{D}_n) \mathbb{P}(\mathcal{P}' \in T_n(\Theta_2), \mathcal{P} \in T_n(\Theta_2) | \mathcal{D}_n)}{\mathbb{P}(\mathcal{P} \in T_n(\Theta_1) | \mathcal{D}_n) \mathbb{P}(\mathcal{P} \in T_n(\Theta_2) | \mathcal{D}_n)} \\
&= \mathbb{P}(\mathcal{P}' \in T_n(\Theta_1) | \mathcal{P} \in T_n(\Theta_1), \mathcal{D}_n) \mathbb{P}(\mathcal{P}' \in T_n(\Theta_2) | \mathcal{P} \in T_n(\Theta_2), \mathcal{D}_n) \\
&= \mathbb{P}(\mathcal{P}' \in T_n(\Theta_1) | \mathcal{P} \in T_n(\Theta_1), \mathcal{D}_n)^2
\end{aligned}$$

we can rewrite  $C_M$

$$\begin{aligned}
C_M &= \mathbb{P}(\mathcal{P}' \in T_n(\Theta_1) | \mathcal{P} \in T_n(\Theta_1), \mathcal{D}_n)^2 \times \Delta_{M,1} \\
&\quad + 2\mathbb{P}(\mathcal{P}' \in T_n(\Theta_1) | \mathcal{P} \in T_n(\Theta_1), \mathcal{D}_n) \\
&\quad \quad \times \mathbb{P}(\mathcal{P}' \in T_n(\Theta_1) | \mathcal{P} \notin T_n(\Theta_1), \mathcal{D}_n) \times \Delta_{M,2} \\
&\quad + \mathbb{P}(\mathcal{P}' \in T_n(\Theta_1) | \mathcal{P} \notin T_n(\Theta_1), \mathcal{D}_n)^2 \times \Delta_{M,3},
\end{aligned}$$

where

$$\begin{aligned}
\Delta_{M,1} &= \mathbb{P}(\mathcal{P} \in T_n(\Theta_1), \mathcal{P} \in T_n(\Theta_2) | \hat{p}_{M,n}(\mathcal{P}) \leq p_0, \mathcal{D}_n) \\
&\quad - \mathbb{P}(\mathcal{P} \in T_n(\Theta_1) | \hat{p}_{M,n}(\mathcal{P}) \leq p_0, \mathcal{D}_n)^2, \\
\Delta_{M,2} &= \mathbb{P}(\mathcal{P} \in T_n(\Theta_1), \mathcal{P} \notin T_n(\Theta_2) | \hat{p}_{M,n}(\mathcal{P}) \leq p_0, \mathcal{D}_n) \\
&\quad - \mathbb{P}(\mathcal{P} \in T_n(\Theta_1) | \hat{p}_{M,n}(\mathcal{P}) \leq p_0, \mathcal{D}_n) \\
&\quad \quad (1 - \mathbb{P}(\mathcal{P} \in T_n(\Theta_1) | \hat{p}_{M,n}(\mathcal{P}) \leq p_0, \mathcal{D}_n)), \\
\Delta_{M,3} &= \mathbb{P}(\mathcal{P} \notin T_n(\Theta_1), \mathcal{P} \notin T_n(\Theta_2) | \hat{p}_{M,n}(\mathcal{P}) \leq p_0, \mathcal{D}_n) \\
&\quad - \mathbb{P}(\mathcal{P} \notin T_n(\Theta_1) | \hat{p}_{M,n}(\mathcal{P}) \leq p_0, \mathcal{D}_n)^2.
\end{aligned}$$

We first consider the term

$$\begin{aligned}
\Delta_{M,1} &= \mathbb{P}(\mathcal{P} \in T_n(\Theta_1), \mathcal{P} \in T_n(\Theta_2) | \hat{p}_{M,n}(\mathcal{P}) \leq p_0, \mathcal{D}_n) \\
&\quad - \mathbb{P}(\mathcal{P} \in T_n(\Theta_1) | \hat{p}_{M,n}(\mathcal{P}) \leq p_0, \mathcal{D}_n)^2
\end{aligned}$$

Equation (D.13) directly gives,

$$\lim_{M \rightarrow \infty} \mathbb{P}(\mathcal{P} \in T_n(\Theta_1) | \hat{p}_{M,n}(\mathcal{P}) \leq p_0, \mathcal{D}_n)^2 = p_0^2. \quad (\text{D.15})$$

On the other hand

$$\begin{aligned}
& \mathbb{P}(\mathcal{P} \in T_n(\Theta_1), \mathcal{P} \in T_n(\Theta_2) | \hat{p}_{M,n}(\mathcal{P}) \leq p_0, \mathcal{D}_n) \\
&= p_n(\mathcal{P})^2 \frac{\mathbb{P}(\hat{p}_{M,n}(\mathcal{P}) \leq p_0 | \mathcal{P} \in T_n(\Theta_1), \mathcal{P} \in T_n(\Theta_2), \mathcal{D}_n)}{\mathbb{P}(\hat{p}_{M,n}(\mathcal{P}) \leq p_0 | \mathcal{D}_n)} \\
&= p_n(\mathcal{P})^2 \frac{\Phi(Mp_0 - 2, M - 2, p_n(\mathcal{P}))}{\Phi(Mp_0, M, p_n(\mathcal{P}))}.
\end{aligned}$$

Again, as for equation (D.11), we can express the cdf of the binomial distribution using the hypergeometric function  $F$

$$\begin{aligned} & \mathbb{P}(\mathcal{P} \in T_n(\Theta_1), \mathcal{P} \in T_n(\Theta_2) | \hat{p}_{M,n}(\mathcal{P}) \leq p_0, \mathcal{D}_n) \\ &= p_0^2 \left( 1 + \frac{p_0 - 1}{p_0(M-1)} \right) \frac{F(M-1, 1, M(1-p_0)+1, 1-p_n(\mathcal{P}))}{F(M+1, 1, M(1-p_0)+1, 1-p_n(\mathcal{P}))}, \end{aligned} \quad (\text{D.16})$$

and from Lemma 10,

$$\lim_{M \rightarrow \infty} \frac{F(M-1, 1, M(1-p_0)+1, 1-p_n(\mathcal{P}))}{F(M+1, 1, M(1-p_0)+1, 1-p_n(\mathcal{P}))} = 1,$$

that is

$$\lim_{M \rightarrow \infty} \mathbb{P}(\mathcal{P} \in T_n(\Theta_1), \mathcal{P} \in T_n(\Theta_2) | \hat{p}_{M,n}(\mathcal{P}) \leq p_0, \mathcal{D}_n) = p_0^2. \quad (\text{D.17})$$

Using equations (D.15) and (D.17), we conclude

$$\lim_{M \rightarrow \infty} \Delta_{M,1} = 0.$$

We follow the same reasoning for  $\Delta_{M,3}$ , equation (D.13) gives

$$\lim_{M \rightarrow \infty} \mathbb{P}(\mathcal{P} \notin T_n(\Theta_1) | \hat{p}_{M,n}(\mathcal{P}) \leq p_0, \mathcal{D}_n)^2 = (1-p_0)^2. \quad (\text{D.18})$$

On the other hand,

$$\begin{aligned} & \mathbb{P}(\mathcal{P} \notin T_n(\Theta_1), \mathcal{P} \notin T_n(\Theta_2) | \hat{p}_{M,n}(\mathcal{P}) \leq p_0, \mathcal{D}_n) \\ &= (1-p_0)^2 \left( 1 - \frac{p_0}{M-1} \right) \frac{F(M-1, 1, M(1-p_0)-11, 1-p_n(\mathcal{P}))}{F(M+1, 1, M(1-p_0)+1, 1-p_n(\mathcal{P}))} \end{aligned}$$

From Lemma 10,

$$\lim_{M \rightarrow \infty} \mathbb{P}(\mathcal{P} \notin T_n(\Theta_1), \mathcal{P} \notin T_n(\Theta_2) | \hat{p}_{M,n}(\mathcal{P}) \leq p_0, \mathcal{D}_n) = (1-p_0)^2 \quad (\text{D.19})$$

And finally  $\lim_{M \rightarrow \infty} \Delta_{M,3} = 0$ . The term  $\Delta_{M,2}$  can be treated in a similar way, since equation (D.13) gives

$$\begin{aligned} & \lim_{M \rightarrow \infty} \mathbb{P}(\mathcal{P} \in T_n(\Theta_1) | \hat{p}_{M,n}(\mathcal{P}) \leq p_0, \mathcal{D}_n) \mathbb{P}(\mathcal{P} \notin T_n(\Theta_1) | \hat{p}_{M,n}(\mathcal{P}) \leq p_0, \mathcal{D}_n) \\ &= p_0(1-p_0). \end{aligned}$$

Simple identity shows

$$\begin{aligned} & \mathbb{P}(\mathcal{P} \in T_n(\Theta_1), \mathcal{P} \notin T_n(\Theta_2) | \hat{p}_{M,n}(\mathcal{P}) \leq p_0, \mathcal{D}_n) \\ &= \frac{1}{2} \left( 1 - \mathbb{P}(\mathcal{P} \notin T_n(\Theta_1), \mathcal{P} \notin T_n(\Theta_2) | \hat{p}_{M,n}(\mathcal{P}) \leq p_0, \mathcal{D}_n) \right. \\ & \quad \left. - \mathbb{P}(\mathcal{P} \in T_n(\Theta_1), \mathcal{P} \in T_n(\Theta_2) | \hat{p}_{M,n}(\mathcal{P}) \leq p_0, \mathcal{D}_n) \right). \end{aligned}$$

Taking the limit of the previous equation and using equations (D.17) and (D.19), we get

$$\begin{aligned} \lim_{M \rightarrow \infty} \mathbb{P}(\mathcal{P} \in T_n(\Theta_1), \mathcal{P} \notin T_n(\Theta_2) | \hat{p}_{M,n}(\mathcal{P}) \leq p_0, \mathcal{D}_n) \\ = p_0(1 - p_0). \end{aligned} \quad (\text{D.20})$$

Using (D.13) and (D.20),  $\lim_{M \rightarrow \infty} \Delta_{M,2} = 0$ . Since  $\Delta_{M,1}, \Delta_{M,2}, \Delta_{M,3} \rightarrow 0$ , we obtain  $\lim_{M \rightarrow \infty} C_M = 0$ , that is,

$$\lim_{M \rightarrow \infty} \mathbb{V}[\hat{p}_{M,n}(\mathcal{P}') | \hat{p}_{M,n}(\mathcal{P}) \leq p_0, \mathcal{D}_n] = 0. \quad (\text{D.21})$$

Finally combining equations (D.14) and (D.21),

$$\begin{aligned} \lim_{M \rightarrow \infty} \mathbb{P}(\hat{p}_{M,n}(\mathcal{P}') > p_0 | \hat{p}_{M,n}(\mathcal{P}) \leq p_0, \mathcal{D}_n) \\ = \mathbb{1}_{p_n(\mathcal{P}') + \rho_n(\mathcal{P}, \mathcal{P}') \frac{\sigma_n(\mathcal{P}')}{\sigma_n(\mathcal{P})} (p_0 - p_n(\mathcal{P})) > p_0} \end{aligned}$$

**Case 2:**  $p_n(\mathcal{P}) \leq p_0$  By the law of large numbers,  $\lim_{M \rightarrow \infty} \hat{p}_{M,n}(\mathcal{P}) = p_n(\mathcal{P})$  in probability, and consequently  $\lim_{M \rightarrow \infty} \mathbb{P}(\hat{p}_{M,n}(\mathcal{P}) \leq p_0) = 1$ . Additionally, we can simply write

$$\begin{aligned} \mathbb{P}(\hat{p}_{M,n}(\mathcal{P}') > p_0 | \hat{p}_{M,n}(\mathcal{P}) \leq p_0, \mathcal{D}_n) \\ = \frac{\mathbb{P}(\hat{p}_{M,n}(\mathcal{P}') > p_0, \hat{p}_{M,n}(\mathcal{P}) \leq p_0 | \mathcal{D}_n)}{\mathbb{P}(\hat{p}_{M,n}(\mathcal{P}) \leq p_0, \mathcal{D}_n)} \end{aligned}$$

Again, by the law of large numbers,  $\lim_{M \rightarrow \infty} \hat{p}_{M,n}(\mathcal{P}') = p_n(\mathcal{P}')$  in probability. Then, if  $p_n(\mathcal{P}') > p_0$ ,  $\lim_{M \rightarrow \infty} \mathbb{P}(\hat{p}_{M,n}(\mathcal{P}') > p_0) = 1$ , and it follows that  $\lim_{M \rightarrow \infty} \mathbb{P}(\hat{p}_{M,n}(\mathcal{P}') > p_0, \hat{p}_{M,n}(\mathcal{P}) \leq p_0 | \mathcal{D}_n) = 1$ . If  $p_n(\mathcal{P}') \leq p_0$ ,  $\lim_{M \rightarrow \infty} \mathbb{P}(\hat{p}_{M,n}(\mathcal{P}') > p_0) = 0$ , and consequently  $\lim_{M \rightarrow \infty} \mathbb{P}(\hat{p}_{M,n}(\mathcal{P}') > p_0, \hat{p}_{M,n}(\mathcal{P}) \leq p_0 | \mathcal{D}_n) = 0$ . This can be compacted under the form

$$\lim_{M \rightarrow \infty} \mathbb{P}(\hat{p}_{M,n}(\mathcal{P}') > p_0 | \hat{p}_{M,n}(\mathcal{P}) \leq p_0, \mathcal{D}_n) = \mathbb{1}_{p_n(\mathcal{P}') > p_0}.$$

The proof for the case  $\mathbb{P}[\hat{p}_{M,n}(\mathcal{P}') > p_0 | \hat{p}_{M,n}(\mathcal{P}) > p_0, \mathcal{D}_n]$  is similar.  $\square$