



HAL
open science

Consistency of Random Forests

Erwan Scornet, Gérard Biau, Jean-Philippe Vert

► **To cite this version:**

Erwan Scornet, Gérard Biau, Jean-Philippe Vert. Consistency of Random Forests. 2014. hal-00990008v3

HAL Id: hal-00990008

<https://hal.science/hal-00990008v3>

Preprint submitted on 21 May 2015 (v3), last revised 7 Aug 2015 (v4)

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.

Consistency of Random Forests

Erwan Scornet

Sorbonne Universités, UPMC Univ Paris 06, F-75005, Paris, France
erwan.scornet@upmc.fr

G erard Biau

Sorbonne Universit es, UPMC Univ Paris 06, F-75005, Paris, France
  Institut universitaire de France
gerard.biau@upmc.fr

Jean-Philippe Vert

*MINES ParisTech, PSL-Research University, CBIO-Centre for
Computational Biology, F-77300, Fontainebleau, France*
  Institut Curie, Paris, F-75248, France
  U900, INSERM, Paris, F-75248, France
jean-philippe.vert@mines-paristech.fr

Abstract

Random forests are a learning algorithm proposed by [Breiman \(2001\)](#) which combines several randomized decision trees and aggregates their predictions by averaging. Despite its wide usage and outstanding practical performance, little is known about the mathematical properties of the procedure. This disparity between theory and practice originates in the difficulty to simultaneously analyze both the randomization process and the highly data-dependent tree structure. In the present paper, we take a step forward in forest exploration by proving a consistency result for Breiman’s [\(2001\)](#) original algorithm in the context of additive regression models. Our analysis also sheds an interesting light on how random forests can nicely adapt to sparsity in high-dimensional settings.

Index Terms — Random forests, randomization, consistency, additive model, sparsity, dimension reduction.

2010 Mathematics Subject Classification: 62G05, 62G20.

1 Introduction

Random forests are an ensemble learning method for classification and regression that constructs a number of randomized decision trees during the

training phase and predicts by averaging the results. Since its publication in the seminal paper of Breiman (2001), the procedure has become a major data analysis tool, that performs well in practice in comparison with many standard methods. What has greatly contributed to the popularity of forests is the fact that they can be applied to a wide range of prediction problems and have few parameters to tune. Aside from being simple to use, the method is generally recognized for its accuracy and its ability to deal with small sample sizes, high-dimensional feature spaces, and complex data structures. The random forest methodology has been successfully involved in many practical problems, including air quality prediction (winning code of the EMC data science global hackathon in 2012, see <http://www.kaggle.com/c/dsg-hackathon>), chemoinformatics (Svetnik et al., 2003), ecology (Prasad et al., 2006; Cutler et al., 2007), 3D object recognition (Shotton et al., 2011), and bioinformatics (Díaz-Uriarte and de Andrés, 2006), just to name a few. In addition, many variations on the original algorithm have been proposed to improve the calculation time while maintaining good prediction accuracy (see, e.g., Geurts et al., 2006; Amaratunga et al., 2008). Breiman’s forests have also been extended to quantile estimation (Meinshausen, 2006), survival analysis (Ishwaran et al., 2008), and ranking prediction (Cléménçon et al., 2013).

On the theoretical side, the story is less conclusive and, regardless of their extensive use in practical settings, little is known about the mathematical properties of random forests. To date, most studies have concentrated on isolated parts or simplified versions of the procedure. The most celebrated theoretical result is that of Breiman (2001), which offers an upper bound on the generalization error of forests in terms of correlation and strength of the individual trees. This was followed by a technical note (Breiman, 2004), that focuses on a stylized version of the original algorithm. A critical step was subsequently taken by Lin and Jeon (2006), who established an interesting connection between random forests and a particular class of nearest neighbor predictors, further explored by Biau and Devroye (2010). In recent years, various theoretical studies (e.g., Biau et al., 2008; Ishwaran and Kogalur, 2010; Biau, 2012; Genuer, 2012; Zhu et al., 2012) have been performed, analyzing consistency of simplified models, and moving ever closer to practice. Recent attempts towards narrowing the gap between theory and practice are by Denil et al. (2013), who proves the first consistency result for online random forests, and by Wager (2014) and Mentch and Hooker (2014) who study some asymptotic properties of forests.

The difficulty to properly analyze random forests can be explained by the black-box nature of the procedure, which is actually a subtle combination of

different components it is illusory to analyze separately. Among the forest essential ingredients, both bagging (Breiman, 1996) and the Classification And Regression Trees (CART)-split criterion (Breiman et al., 1984) play a critical role. Bagging (a contraction of bootstrap-aggregating) is a general aggregation scheme which proceeds by generating subsamples from the original data set, constructing a predictor from each resample and deciding by averaging. It is one of the most effective computationally intensive procedures to improve on unstable estimates, especially for large, high-dimensional data sets where finding a good model in one step is impossible because of the complexity and scale of the problem (Bühlmann and Yu, 2002; Kleiner et al., 2012; Wager et al., 2013). On the other hand, the CART-split selection, originated from the most influential CART algorithm of Breiman et al. (1984), is used in the construction of the individual trees to choose the best cuts perpendicular to the axes. At each node of each tree, the best cut is selected by optimizing the CART-split criterion, based on the notion of Gini impurity (classification) and prediction squared error (regression).

Yet, while bagging and the CART-splitting scheme play a key role in the random forest mechanism, both are difficult to analyze, thereby explaining why theoretical studies have considered so far simplified versions of the original procedure. This is often done by simply ignoring the bagging step and by replacing the CART-split selection by a more elementary cut protocol. Besides, in Breiman’s forests, each leaf (that is, a terminal node) of the individual trees contains a fixed pre-specified number of observations (this parameter, called `nodesize` in the R package `randomForests`, is usually chosen between 1 and 5). There is also an extra parameter in the algorithm which allows to control the total number of leaves (this parameter is called `maxnode` in the R package and has, by default, no effect on the procedure). The combination of these various components makes the algorithm difficult to analyze with rigorous mathematics. As a matter of fact, most authors focus on simplified, data-independent procedures, thus creating a gap between theory and practice.

Motivated by the above discussion, we study in the present paper some asymptotic properties of Breiman’s (2001) algorithm in the context of additive regression models. We prove the \mathbb{L}^2 consistency of random forests, which gives a first basic theoretical guarantee of efficiency for this algorithm. Up to our knowledge, this is the first consistency result for Breiman’s (2001) original procedure. Our approach rests upon a detailed analysis of the behavior of the cells generated by CART-split selection as the sample size grows. It turns out that a good control of the regression function variation inside each cell, together with a proper choice of the total number of leaves (Theorem

3.1) or a proper choice of the subsampling rate (Theorem 3.2) are sufficient to ensure the forest consistency in a \mathbb{L}^2 sense. Also, our analysis shows that random forests can adapt to a sparse framework, when the ambient dimension p is large but only a smaller number of coordinates carry out information.

The paper is organized as follows. In Section 2, we introduce some notations and describe the random forest method. The main asymptotic results are presented in Section 3 and further discussed in Section 4. Section 5 is devoted to the main proofs, and technical results are postponed to Section 6.

2 Random forests

The general framework is \mathbb{L}^2 regression estimation, in which an input random vector $\mathbf{X} \in [0, 1]^p$ is observed, and the goal is to predict the square integrable random response $Y \in \mathbb{R}$ by estimating the regression function $m(\mathbf{x}) = \mathbb{E}[Y|\mathbf{X} = \mathbf{x}]$. To this aim, we assume given a training sample $\mathcal{D}_n = (\mathbf{X}_1, Y_1), \dots, (\mathbf{X}_n, Y_n)$ of $[0, 1]^p \times \mathbb{R}$ -valued independent random variables distributed as the independent prototype pair (\mathbf{X}, Y) . The objective is to use the data set \mathcal{D}_n to construct an estimate $m_n : [0, 1]^p \rightarrow \mathbb{R}$ of the function m . In this respect, we say that a regression function estimate m_n is \mathbb{L}^2 consistent if $\mathbb{E}[m_n(\mathbf{X}) - m(\mathbf{X})]^2 \rightarrow 0$ as $n \rightarrow \infty$ (where the expectation is over \mathbf{X} and \mathcal{D}_n).

A random forest is a predictor consisting of a collection of M randomized regression trees. For the j -th tree in the family, the predicted value at the query point \mathbf{x} is denoted by $m_n(\mathbf{x}; \Theta_j, \mathcal{D}_n)$, where $\Theta_1, \dots, \Theta_M$ are independent random variables, distributed as a generic random variable Θ and independent of \mathcal{D}_n . In practice, this variable is used to resample the training set prior to the growing of individual trees and to select the successive candidate directions for splitting. The trees are combined to form the (finite) forest estimate

$$m_{M,n}(\mathbf{x}; \Theta_1, \dots, \Theta_M, \mathcal{D}_n) = \frac{1}{M} \sum_{j=1}^M m_n(\mathbf{x}; \Theta_j, \mathcal{D}_n). \quad (1)$$

Since in practice we can choose M as large as possible, we study in this paper the property of the infinite forest estimate obtained as the limit of (1) when the number of trees M grows to infinity:

$$m_n(\mathbf{x}; \mathcal{D}_n) = \mathbb{E}_{\Theta} [m_n(\mathbf{x}; \Theta, \mathcal{D}_n)],$$

where \mathbb{E}_Θ denotes expectation with respect to the random parameter Θ , conditional on \mathcal{D}_n . This operation is justified by the law of large numbers, which asserts that, almost surely, conditional on \mathcal{D}_n ,

$$\lim_{M \rightarrow \infty} m_{n,M}(\mathbf{x}; \Theta_1, \dots, \Theta_M, \mathcal{D}_n) = m_n(\mathbf{x}; \mathcal{D}_n)$$

(see, e.g., [Scornet, 2014](#), for details). In the sequel, to lighten notation, we will simply write $m_n(\mathbf{x})$ instead of $m_n(\mathbf{x}; \mathcal{D}_n)$.

In Breiman's (2001) original forests, each node of a single tree is associated with a hyper-rectangular cell. At each step of the tree construction, the collection of cells forms a partition of $[0, 1]^p$. The root of the tree is $[0, 1]^p$ itself, and each tree is grown as explained in **Algorithm 1**.

This algorithm has three parameters:

1. $m_{\text{try}} \in \{1, \dots, p\}$, which is the number of pre-selected directions for splitting;
2. $a_n \in \{1, \dots, n\}$, which is the number of sampled data points in each tree;
3. $t_n \in \{1, \dots, a_n\}$, which is the number of leaves in each tree.

By default, in the original procedure, the parameter m_{try} is set to $p/3$, a_n is set to n (resampling is done with replacement), and $t_n = a_n$. However, in our approach, resampling is done without replacement and the parameters a_n and t_n can be different from their default values.

In a word, the algorithm works by growing M different trees as follows. For each tree, a_n data points are drawn at random without replacement from the original data set; then, at each cell of every tree, a split is chosen by maximizing the CART-criterion (see below); finally, the construction of every tree is stopped when the total number of cells in the tree reaches the value t_n (therefore, each cell contains exactly one point in the case $t_n = a_n$).

Algorithm 1: Breiman’s random forest predicted value at \mathbf{x} .

Input: Training set \mathcal{D}_n , number of trees $M > 0$, $m_{\text{try}} \in \{1, \dots, p\}$,
 $a_n \in \{1, \dots, n\}$, $t_n \in \{1, \dots, a_n\}$, and $\mathbf{x} \in [0, 1]^p$.

Output: Prediction of the random forest at \mathbf{x} .

```

1 for  $j = 1, \dots, M$  do
2   Select  $a_n$  points, without replacement, uniformly in  $\mathcal{D}_n$ .
3   Set  $\mathcal{P}_0 = \{[0, 1]^p\}$  the partition associated with the root of the tree.
4   Set  $n_{\text{nodes}} = 1$  and  $\text{level} = 0$ .
5   while  $n_{\text{nodes}} < t_n$  do
6     if  $\mathcal{P}_{\text{level}} = \emptyset$  then
7       |  $\text{level} = \text{level} + 1$ 
8     else
9       | Let  $A$  be the first element in  $\mathcal{P}_{\text{level}}$ .
10      | if  $A$  contains exactly one point then
11        | |  $\mathcal{P}_{\text{level}} \leftarrow \mathcal{P}_{\text{level}} \setminus \{A\}$ ;  $\mathcal{P}_{\text{level}+1} \leftarrow \mathcal{P}_{\text{level}+1} \cup \{A\}$ 
12      | else
13        | | Select uniformly, without replacement, a subset
14        | |  $\mathcal{M}_{\text{try}} \subset \{1, \dots, p\}$  of cardinality  $m_{\text{try}}$ .
15        | | Select the best split in  $A$  by optimizing the CART-split
16        | | criterion along the coordinates in  $\mathcal{M}_{\text{try}}$  (see details below).
17        | | Cut the cell  $A$  according to the best split. Call  $A_L$  and  $A_R$ 
18        | | the two resulting cell.
19        | |  $\mathcal{P}_{\text{level}} \leftarrow \mathcal{P}_{\text{level}} \setminus \{A\}$ 
20        | |  $\mathcal{P}_{\text{level}+1} \leftarrow \mathcal{P}_{\text{level}+1} \cup \{A_L\} \cup \{A_R\}$ 
21        | |  $n_{\text{nodes}} = n_{\text{nodes}} + 1$ 
22      | end
23    end
24  end
25  Compute the predicted value  $m_n(\mathbf{x}; \Theta_j, \mathcal{D}_n)$  at  $\mathbf{x}$  equal to the average of
26  the  $Y_i$ ’s falling in the cell of  $\mathbf{x}$  in partition  $\mathcal{P}_{\text{level}} \cup \mathcal{P}_{\text{level}+1}$ .
27 end
28 Compute the random forest estimate  $m_{M,n}(\mathbf{x}; \Theta_1, \dots, \Theta_M, \mathcal{D}_n)$  at the query
29 point  $\mathbf{x}$  according to (1).

```

We note that the resampling step in **Algorithm 1** (line 2) is done by choosing a_n out of n points (with $a_n \leq n$) without replacement. This is slightly different from the original algorithm, where resampling is done by bootstrapping, that is by choosing n out of n data points with replacement.

Selecting the points “without replacement” instead of “with replacement”

is harmless—in fact, it is just a means to avoid mathematical difficulties induced by the bootstrap.

On the other hand, letting the parameters a_n and t_n depend upon n offers several degrees of freedom which opens the route for establishing consistency of the method. To be precise, we will study in Section 3 the random forest algorithm in two different regimes. The first regime is when $t_n < a_n$, which means that trees are not fully developed. In that case, a proper tuning of t_n ensures the forest’s consistency (Theorem 3.1). The second regime occurs when $t_n = a_n$, i.e. when trees are fully grown. In that case, consistency results from an appropriate choice of the subsample rate a_n/n (Theorem 3.2).

So far, we have not made explicit the CART-split criterion used in **Algorithm 1**. To properly define it, we let A be a generic cell and $N_n(A)$ be the number of data points falling in A . A cut in A is a pair (j, z) , where j is a dimension in $\{1, \dots, p\}$ and z is the position of the cut along the j -th coordinate, within the limits of A . We let \mathcal{C}_A be the set of all such possible cuts in A . Then, with the notation $\mathbf{X}_i = (\mathbf{X}_i^{(1)}, \dots, \mathbf{X}_i^{(p)})$, for any $(j, z) \in \mathcal{C}_A$, the CART-split criterion (Breiman et al., 1984) takes the form

$$L_n(j, z) = \frac{1}{N_n(A)} \sum_{i=1}^n (Y_i - \bar{Y}_A)^2 \mathbf{1}_{\mathbf{x}_i \in A} - \frac{1}{N_n(A)} \sum_{i=1}^n (Y_i - \bar{Y}_{A_L} \mathbf{1}_{\mathbf{x}_i^{(j)} < z} - \bar{Y}_{A_R} \mathbf{1}_{\mathbf{x}_i^{(j)} \geq z})^2 \mathbf{1}_{\mathbf{x}_i \in A}, \quad (2)$$

where $A_L = \{\mathbf{x} \in A : \mathbf{x}^{(j)} < z\}$, $A_R = \{\mathbf{x} \in A : \mathbf{x}^{(j)} \geq z\}$, and \bar{Y}_A (resp., \bar{Y}_{A_L} , \bar{Y}_{A_R}) is the average of the Y_i ’s belonging to A (resp., A_L , A_R), with the convention $0/0 = 0$. At each cell A , the best cut (j_n^*, z_n^*) is finally selected by maximizing $L_n(j, z)$ over \mathcal{M}_{try} and \mathcal{C}_A , that is

$$(j_n^*, z_n^*) \in \arg \max_{\substack{j \in \mathcal{M}_{\text{try}} \\ (j, z) \in \mathcal{C}_A}} L_n(j, z).$$

To remove ties in the argmax, the best cut is always performed along the best cut direction j_n^* , at the middle of two consecutive data points.

3 Main results

We consider an additive regression model satisfying the following properties:

(H1) *The response Y follows*

$$Y = \sum_{j=1}^p m_j(\mathbf{X}^{(j)}) + \varepsilon,$$

where $\mathbf{X} = (\mathbf{X}^{(1)}, \dots, \mathbf{X}^{(p)})$ is uniformly distributed over $[0, 1]^p$, ε is an independent centered Gaussian noise with finite variance $\sigma^2 > 0$, and each component m_j is continuous.

Additive regression models, which extend linear models, were popularized by [Stone \(1985\)](#) and [Hastie and Tibshirani \(1986\)](#). These models, which decompose the regression function as a sum of univariate functions, are flexible and easy to interpret. They are acknowledged for providing a good trade-off between model complexity and calculation time, and were accordingly extensively studied for the last thirty years. Additive models also play an important role in the context of high-dimensional data analysis and sparse modelling, where they are successfully involved in procedures such as the Lasso and various aggregation schemes (for an overview, see, e.g., [Hastie et al., 2009](#)).

Our first result assumes that the total number of leaves t_n in each tree tends to infinity slower than the number of selected data points a_n .

Theorem 3.1. *Assume that **(H1)** is satisfied. Then, provided $a_n \rightarrow \infty$ and $t_n(\log a_n)^9/a_n \rightarrow 0$, random forests are consistent, i.e.,*

$$\lim_{n \rightarrow \infty} \mathbb{E} [m_n(\mathbf{X}) - m(\mathbf{X})]^2 = 0.$$

Noteworthy, [Theorem 3.1](#) still holds with $a_n = n$. In that case, the subsampling step plays no role in the consistency of the method. Indeed, controlling the depth of the trees via the parameter t_n is sufficient to bound the forest error. We note in passing that an easy adaptation of [Theorem 3.1](#) shows that the CART algorithm is consistent under the same assumptions.

The term $(\log a_n)^9$ originates from the Gaussian noise and allows to control the noise tail. In the easier situation where the Gaussian noise is replaced by a bounded random variable, it is easy to see that the term $(\log a_n)^9$ turns into $\log a_n$, a term which accounts for the complexity of the tree partition.

Let us now examine the forest behavior in the second regime, where $t_n = a_n$ (i.e., trees are fully grown) and, as before, subsampling is done at the rate a_n/n . The analysis of this regime turns out to be more complicated, and

rests upon assumption **(H2)** below. We denote by $Z_i = \mathbb{1}_{\mathbf{X} \leftrightarrow \mathbf{X}_i}$ the indicator that \mathbf{X}_i falls in the same cell as \mathbf{X} in the random tree designed with \mathcal{D}_n and the random parameter Θ . Similarly, we let $Z'_j = \mathbb{1}_{\mathbf{X} \leftrightarrow \mathbf{X}_j}$, where Θ' is an independent copy of Θ . Accordingly, we define

$$\begin{aligned} \psi_{i,j}(Y_i, Y_j) &= \mathbb{E} \left[Z_i Z'_j \mid \mathbf{X}, \Theta, \Theta', \mathbf{X}_1, \dots, \mathbf{X}_n, Y_i, Y_j \right] \\ \text{and } \psi_{i,j} &= \mathbb{E} \left[Z_i Z'_j \mid \mathbf{X}, \Theta, \Theta', \mathbf{X}_1, \dots, \mathbf{X}_n \right]. \end{aligned}$$

Finally, for any random variables W_1, W_2, Z , we denote by $\text{Corr}(W_1, W_2 \mid Z)$ the conditional correlation coefficient (whenever it exists).

(H2) Let $Z_{i,j} = (Z_i, Z'_j)$. Then, one of the following two conditions holds:

(H2.1) One has

$$\lim_{n \rightarrow \infty} (\log a_n)^{2d-2} (\log n)^2 \mathbb{E} \left[\max_{\substack{i,j \\ i \neq j}} |\psi_{i,j}(Y_i, Y_j) - \psi_{i,j}| \right]^2 = 0.$$

(H2.2) There exist a constant $C > 0$ and a sequence $(\gamma_n)_n \rightarrow 0$ such that, almost surely,

$$\max_{\ell_1, \ell_2=0,1} \frac{|\text{Corr}(Y_i - m(\mathbf{X}_i), \mathbb{1}_{Z_{i,j}=(\ell_1, \ell_2)} \mid \mathbf{X}_i, \mathbf{X}_j, Y_j)|}{\mathbb{P}^{1/2}[Z_{i,j} = (\ell_1, \ell_2) \mid \mathbf{X}_i, \mathbf{X}_j, Y_j]} \leq \gamma_n,$$

and

$$\max_{\ell_1=0,1} \frac{|\text{Corr}((Y_i - m(\mathbf{X}_i))^2, \mathbb{1}_{Z_i=\ell_1} \mid \mathbf{X}_i)|}{\mathbb{P}^{1/2}[Z_i = \ell_1 \mid \mathbf{X}_i]} \leq C.$$

Despite their technical aspect, statements **(H2.1)** and **(H2.2)** have simple interpretations. To understand the meaning of **(H2.1)**, let us replace the Gaussian noise by a bounded random variable. A close inspection of Lemma 4 shows that **(H2.1)** may be simply replaced by

$$\lim_{n \rightarrow \infty} \mathbb{E} \left[\max_{\substack{i,j \\ i \neq j}} |\psi_{i,j}(Y_i, Y_j) - \psi_{i,j}| \right]^2 = 0.$$

Therefore, **(H2.1)** means that the influence of two Y-values on the probability of connection of two couples of random points tends to zero as $n \rightarrow \infty$.

As for assumption **(H2.2)**, it holds whenever the correlation between the noise and the probability of connection of two couples of random points vanishes fast enough, as $n \rightarrow \infty$. Note that, in the simple case where the partition is independent of the Y_i 's, the correlations in **(H2.2)** are zero, so that **(H2)** is trivially satisfied. It is also verified in the noiseless case, that is, when $Y = m(\mathbf{X})$. However, in the most general context, the partitions strongly depend on the whole sample \mathcal{D}_n and, unfortunately, we do not know whether **(H2)** is satisfied or not.

Theorem 3.2. *Assume that **(H1)** and **(H2)** are satisfied and let $t_n = a_n$. Then, provided $a_n \rightarrow \infty$ and $a_n \log n/n \rightarrow 0$, random forests are consistent, i.e.,*

$$\lim_{n \rightarrow \infty} \mathbb{E} [m_n(\mathbf{X}) - m(\mathbf{X})]^2 = 0.$$

Up to our knowledge, apart from the fact that bootstrapping is replaced by subsampling, Theorem 3.1 and Theorem 3.2 are the first consistency results for Breiman's (2001) forests. Indeed, most models studied so far are designed independently of \mathcal{D}_n and are, consequently, an unrealistic representation of the true procedure. In fact, understanding Breiman's random forest behavior deserves a more involved mathematical treatment. Section 4 below offers a thorough description of the various mathematical forces in action.

Our study also sheds some interesting light on the behavior of forests when the ambient dimension p is large but the true underlying dimension of the model is small. To see how, assume that the additive model **(H1)** satisfies a sparsity constraint of the form

$$Y = \sum_{j=1}^S m_j(\mathbf{X}^{(j)}) + \varepsilon,$$

where $S < p$ represents the true, but unknown, dimension of the model. Thus, among the p original features, it is assumed that only the first (without loss of generality) S variables are informative. Put differently, Y is assumed to be independent of the last $(p - S)$ variables. In this dimension reduction context, the ambient dimension p can be very large, but we believe that the representation is sparse, i.e., that few components of m are non-zero. As such, the value S characterizes the sparsity of the model: the smaller S , the sparser m .

Proposition 1 below shows that random forests nicely adapt to the sparsity setting by asymptotically performing, with high probability, splits along the S informative variables.

In this proposition, we set $m_{\text{try}} = p$ and, for all k , we denote by $j_{1,n}(\mathbf{X}), \dots, j_{k,n}(\mathbf{X})$ the first k cut directions used to construct the cell containing \mathbf{X} , with the convention that $j_{q,n}(\mathbf{X}) = \infty$ if the cell has been cut strictly less than q times.

Proposition 1. *Assume that (H1) is satisfied. Let $k \in \mathbb{N}^*$ and $\xi > 0$. Assume that there is no interval $[a, b]$ and no $j \in \{1, \dots, S\}$ such that m_j is constant on $[a, b]$. Then, with probability $1 - \xi$, for all n large enough, we have, for all $1 \leq q \leq k$,*

$$j_{q,n}(\mathbf{X}) \in \{1, \dots, S\}.$$

This proposition provides an interesting perspective on why random forests are still able to do a good job in high-dimensional settings. Since the algorithm selects splits mostly along informative variables, everything happens as if data were projected onto the vector space generated by the S informative variables. Therefore, forests are likely to only depend upon these S variables, which supports the fact that they have good performance in sparse framework.

4 Discussion

One of the main difficulties in assessing the mathematical properties of Breiman's (2001) forests is that the construction process of the individual trees strongly depends on both the X_i 's and the Y_i 's. For partitions that are independent of the Y_i 's, consistency can be shown by relatively simple means via Stone's (1977) theorem for local averaging estimates (see also Györfi et al., 2002, Chapter 6). However, our partitions and trees depend upon the Y -values in the data. This makes things complicated, but mathematically interesting too. Thus, logically, the proof of Theorem 3.2 starts with an adaptation of Stone's (1977) theorem tailored for random forests, whereas the proof of Theorem 3.1 is based on consistency results of data-dependent partitions developed by Nobel (1996).

Both theorems rely on Proposition 2 below which stresses an important feature of the random forest mechanism. It states that the variation of the regression function m within a cell of a random tree is small provided n is large enough. To this aim, we define, for any cell A , the variation of m within A as

$$\Delta(m, A) = \sup_{\mathbf{x}, \mathbf{x}' \in A} |m(\mathbf{x}) - m(\mathbf{x}')|.$$

Furthermore, we denote by $A_n(\mathbf{X}, \Theta)$ the cell of a tree built with random parameter Θ that contains the point \mathbf{X} .

Proposition 2. *Assume that (H1) holds. Then, for all $\rho, \xi > 0$, there exists $N \in \mathbb{N}^*$ such that, for all $n > N$,*

$$\mathbb{P} [\Delta(m, A_n(\mathbf{X}, \Theta)) \leq \xi] \geq 1 - \rho.$$

It should be noted that in the standard, Y -independent analysis of partitioning regression function estimates, the variance is controlled by letting the diameters of the tree cells tend to zero in probability. Instead of such a geometrical assumption, Proposition 2 ensures that the variation of m inside a cell is small, thereby forcing the approximation error of the forest to asymptotically approach zero.

While Proposition 2 offers a good control of the approximation error of the forest in both regimes, a separated analysis is required for the estimation error. In regime 1 (Theorem 3.1), the parameter t_n allows to control the structure of the tree. This is in line with standard tree consistency approaches (see, e.g., Chapter 20 in Devroye et al., 1996). Things are different for the second regime (Theorem 3.2), in which individual trees are fully grown. In that case, the estimation error is controlled by forcing the subsampling rate a_n/n to be $o(1/\log n)$, which is a more unusual requirement and deserves some remarks.

At first, we note that the $\log n$ term in Theorem 3.2 is used to control the Gaussian noise ε . Thus, if the noise is assumed to be a bounded random variable, then the $\log n$ term disappears, and the condition reduces to $a_n/n \rightarrow 0$. The requirement $a_n \log n/n \rightarrow 0$ guarantees that every single observation (\mathbf{X}_i, Y_i) is used in the tree construction with a probability that becomes small with n . It also implies that the query point \mathbf{x} is not connected to the same data point in a high proportion of trees. If not, the predicted value at \mathbf{x} would be too much influenced by one single pair (\mathbf{X}_i, Y_i) , making the forest inconsistent. In fact, the proof of Theorem 3.2 reveals that the estimation error of a forest estimate is small as soon as the maximum probability of connection between the query point and all observations is small. Thus, the assumption on the subsampling rate is just a convenient way to control these probabilities, by ensuring that partitions are dissimilar enough (i.e. by ensuring that \mathbf{x} is connected with many data points through the forest). This idea of diversity among trees was introduced by Breiman (2001), but is generally difficult to analyse. In our approach, the subsampling is the key component for imposing tree diversity.

Theorem 3.2 comes at the price of assumption **(H2)**, for which we do not know if it is valid in all generality. On the other hand, Theorem 3.2, which mimics almost perfectly the algorithm used in practice, is an important step towards understanding Breiman’s random forests. Contrary to most previous works, Theorem 3.2 assumes that there is only one observation per leaf of each individual tree. This implies that the single trees are eventually not consistent, since standard conditions for tree consistency require that the number of observations in the terminal nodes tends to infinity as n grows (see, e.g., Devroye et al., 1996; Györfi et al., 2002). Thus, the random forest algorithm aggregates rough individual tree predictors to build a provably consistent general architecture.

It is also interesting to note that our results (in particular Lemma 3) cannot be directly extended to establish the pointwise consistency of random forests, that is, for almost all $\mathbf{x} \in [0, 1]^d$,

$$\lim_{n \rightarrow \infty} \mathbb{E}[m_n(\mathbf{x}) - m(\mathbf{x})]^2 = 0.$$

Fixing $\mathbf{x} \in [0, 1]^d$, the difficulty results from the fact that we do not have a control on the diameter of the cell $A_n(\mathbf{x}, \Theta)$, whereas, since the cells form a partition of $[0, 1]^d$, we have a global control on their diameters. Thus, as highlighted by Wager (2014), random forests can be inconsistent at some fixed point $\mathbf{x} \in [0, 1]^d$, particularly near the edges, while being \mathbb{L}^2 consistent.

Let us finally mention that all results can be extended to the case where ε is a heteroscedastic and sub-Gaussian noise, with for all $\mathbf{x} \in [0, 1]^d$, $\mathbb{V}[\varepsilon | \mathbf{X} = \mathbf{x}] \leq \sigma'^2$, for some constant σ'^2 . All proofs can be readily extended to match this context, at the price of easy technical adaptations.

5 Proof of Theorem 3.1 and Theorem 3.2

For the sake of clarity, proofs of the intermediary results are postponed to Section 6. We start with some notations.

5.1 Notations

In the sequel, to clarify the notations, we will sometimes write $d = (d^{(1)}, d^{(2)})$ to represent a cut (j, z) .

Recall that, for any cell A , \mathcal{C}_A is the set of all possible cuts in A . Thus, with this notation, $\mathcal{C}_{[0,1]^p}$ is just the set of all possible cuts at the root of the

tree, that is, all possible choices $d = (d^{(1)}, d^{(2)})$ with $d^{(1)} \in \{1, \dots, p\}$ and $d^{(2)} \in [0, 1]$.

More generally, for any $\mathbf{x} \in [0, 1]^p$, we call $\mathcal{A}_k(\mathbf{x})$ the collection of all possible $k \geq 1$ consecutive cuts used to build the cell containing \mathbf{x} . Such a cell is obtained after a sequence of cuts $\mathbf{d}_k = (d_1, \dots, d_k)$, where the dependency of \mathbf{d}_k upon \mathbf{x} is understood. Accordingly, for any $\mathbf{d}_k \in \mathcal{A}_k(\mathbf{x})$, we let $A(\mathbf{x}, \mathbf{d}_k)$ be the cell containing \mathbf{x} built with the particular k -tuple of cuts \mathbf{d}_k . The proximity between two elements \mathbf{d}_k and \mathbf{d}'_k in $\mathcal{A}_k(\mathbf{x})$ will be measured via

$$\|\mathbf{d}_k - \mathbf{d}'_k\|_\infty = \sup_{1 \leq j \leq k} \max(|d_j^{(1)} - d_j'^{(1)}|, |d_j^{(2)} - d_j'^{(2)}|).$$

Accordingly, the distance d_∞ between $\mathbf{d}_k \in \mathcal{A}_k(\mathbf{x})$ and any $\mathcal{A} \subset \mathcal{A}_k(\mathbf{x})$ is

$$d_\infty(\mathbf{d}_k, \mathcal{A}) = \inf_{\mathbf{z} \in \mathcal{A}} \|\mathbf{d}_k - \mathbf{z}\|_\infty.$$

Remember that $A_n(\mathbf{X}, \Theta)$ denotes the cell of a tree containing \mathbf{X} and designed with random parameter Θ . Similarly, $A_{k,n}(\mathbf{X}, \Theta)$ is the same cell but where only the first k cuts are performed ($k \in \mathbb{N}^*$ is a parameter to be chosen later). We also denote by $\hat{\mathbf{d}}_{k,n}(\mathbf{X}, \Theta) = (\hat{d}_{1,n}(\mathbf{X}, \Theta), \dots, \hat{d}_{k,n}(\mathbf{X}, \Theta))$ the k cuts used to construct the cell $A_{k,n}(\mathbf{X}, \Theta)$.

Recall that, for any cell A , the empirical criterion used to split A in the random forest algorithm is defined in (2). For any cut $(j, z) \in \mathcal{C}_A$, we denote the following theoretical version of $L_n(\cdot, \cdot)$ by

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

Observe that $L^*(\cdot, \cdot)$ does not depend upon the training set and that, by the strong law of large numbers, $L_n(j, z) \rightarrow L^*(j, z)$ almost surely as $n \rightarrow \infty$ for all cuts $(j, z) \in \mathcal{C}_A$. Therefore, it is natural to define the best theoretical split (j^*, z^*) of the cell A as

$$(j^*, z^*) \in \arg \min_{\substack{(j,z) \in \mathcal{C}_A \\ j \in \mathcal{M}_{\text{try}}}} L^*(j, z).$$

In view of this criterion, we define the theoretical random forest as before, but with consecutive cuts performed by optimizing $L^*(\cdot, \cdot)$ instead of $L_n(\cdot, \cdot)$. We note that this new forest does depend on Θ through \mathcal{M}_{try} , but not on the sample \mathcal{D}_n . In particular, the stopping criterion for dividing cells has to be

changed in the theoretical random forest; instead of stopping when a cell has a single training point, we impose that each tree of the theoretical forest is stopped at a fixed level $k \in \mathbb{N}^*$. We also let $A_k^*(\mathbf{X}, \Theta)$ be a cell of the theoretical random tree at level k , containing \mathbf{X} , designed with randomness Θ , and resulting from the k theoretical cuts $\mathbf{d}_k^*(\mathbf{X}, \Theta) = (d_1^*(\mathbf{X}, \Theta), \dots, d_k^*(\mathbf{X}, \Theta))$. Since there can exist multiple best cuts at, at least, one node, we call $\mathcal{A}_k^*(\mathbf{X}, \Theta)$ the set of all k -tuples $\mathbf{d}_k^*(\mathbf{X}, \Theta)$ of best theoretical cuts used to build $A_k^*(\mathbf{X}, \Theta)$.

We are now equipped to prove Proposition 2. For clarity reasons, the proof has been divided in three steps. Firstly, we study in Lemma 1 the theoretical random forest. Then we prove in Lemma 3 (via Lemma 2), that theoretical and empirical cuts are close to each other. Proposition 2 is finally established as a consequence of Lemma 1 and Lemma 3. Proofs of these lemmas are to be found in Section 6.

5.2 Proof of Proposition 2

We first need a lemma which states that the variation of $m(\mathbf{X})$ within the cell $A_k^*(\mathbf{X}, \Theta)$ where \mathbf{X} falls, as measured by $\Delta(m, A_k^*(\mathbf{X}, \Theta))$, tends to zero.

Lemma 1. *Assume that (H1) is satisfied. Then, for all $\mathbf{x} \in [0, 1]^p$,*

$$\Delta(m, A_k^*(\mathbf{x}, \Theta)) \rightarrow 0, \quad \text{almost surely, as } k \rightarrow \infty.$$

The next step is to show that cuts in theoretical and original forests are close to each other. To this aim, for any $\mathbf{x} \in [0, 1]^p$ and any k -tuple of cuts $\mathbf{d}_k \in \mathcal{A}_k(\mathbf{x})$, we define

$$\begin{aligned} L_{n,k}(\mathbf{x}, \mathbf{d}_k) &= \frac{1}{N_n(A(\mathbf{x}, \mathbf{d}_{k-1}))} \sum_{i=1}^n (Y_i - \bar{Y}_{A(\mathbf{x}, \mathbf{d}_{k-1})})^2 \mathbb{1}_{\mathbf{X}_i \in A(\mathbf{x}, \mathbf{d}_{k-1})} \\ &\quad - \frac{1}{N_n(A(\mathbf{x}, \mathbf{d}_{k-1}))} \sum_{i=1}^n \left(Y_i - \bar{Y}_{A_L(\mathbf{x}, \mathbf{d}_{k-1})} \mathbb{1}_{\mathbf{X}_i^{(d_k^{(1)})} < d_k^{(2)}} \right. \\ &\quad \left. - \bar{Y}_{A_R(\mathbf{x}, \mathbf{d}_{k-1})} \mathbb{1}_{\mathbf{X}_i^{(d_k^{(1)})} \geq d_k^{(2)}} \right)^2 \mathbb{1}_{\mathbf{X}_i \in A(\mathbf{x}, \mathbf{d}_{k-1})}, \end{aligned}$$

where $A_L(\mathbf{x}, \mathbf{d}_{k-1}) = A(\mathbf{x}, \mathbf{d}_{k-1}) \cap \{\mathbf{z} : \mathbf{z}^{(d_k^{(1)})} < d_k^{(2)}\}$ and $A_R(\mathbf{x}, \mathbf{d}_{k-1}) = A(\mathbf{x}, \mathbf{d}_{k-1}) \cap \{\mathbf{z} : \mathbf{z}^{(d_k^{(1)})} \geq d_k^{(2)}\}$, and where we use the convention $0/0 = 0$ when $A(\mathbf{x}, \mathbf{d}_{k-1})$ is empty. Besides, we let $A(\mathbf{x}, \mathbf{d}_0) = [0, 1]^p$ in the previous

equation. The quantity $L_{n,k}(\mathbf{x}, \mathbf{d}_k)$ is nothing but the criterion to maximize in d_k to find the best k -th cut in the cell $A(\mathbf{x}, \mathbf{d}_{k-1})$. Lemma 2 below ensures that $L_{n,k}(\mathbf{x}, \cdot)$ is stochastically equicontinuous, for all $\mathbf{x} \in [0, 1]^p$. To this aim, for all $\xi > 0$, and for all $\mathbf{x} \in [0, 1]^p$, we denote by $\mathcal{A}_{k-1}^\xi(\mathbf{x}) \subset \mathcal{A}_{k-1}(\mathbf{x})$ the set of all $(k-1)$ -tuples \mathbf{d}_{k-1} such that the cell $A(\mathbf{x}, \mathbf{d}_{k-1})$ contains a hypercube of edge length ξ . Moreover, we let $\bar{\mathcal{A}}_k^\xi(\mathbf{x}) = \{\mathbf{d}_k : \mathbf{d}_{k-1} \in \mathcal{A}_{k-1}^\xi(\mathbf{x})\}$ equipped with the norm $\|\mathbf{d}_k\|_\infty$.

Lemma 2. *Assume that (H1) is satisfied. Fix $\mathbf{x} \in [0, 1]^p$, $k \in \mathbb{N}^*$, and let $\xi > 0$. Then $L_{n,k}(\mathbf{x}, \cdot)$ is stochastically equicontinuous on $\bar{\mathcal{A}}_k^\xi(\mathbf{x})$, that is, for all $\alpha, \rho > 0$, there exists $\delta > 0$ such that*

$$\lim_{n \rightarrow \infty} \mathbb{P} \left[\sup_{\substack{\|\mathbf{d}_k - \mathbf{d}'_k\|_\infty \leq \delta \\ \mathbf{d}_k, \mathbf{d}'_k \in \bar{\mathcal{A}}_k^\xi(\mathbf{x})}} |L_{n,k}(\mathbf{x}, \mathbf{d}_k) - L_{n,k}(\mathbf{x}, \mathbf{d}'_k)| > \alpha \right] \leq \rho.$$

Lemma 2 is then used in Lemma 3 to assess the distance between theoretical and empirical cuts.

Lemma 3. *Assume that (H1) is satisfied. Fix $\xi, \rho > 0$ and $k \in \mathbb{N}^*$. Then there exists $N \in \mathbb{N}^*$ such that, for all $n \geq N$,*

$$\mathbb{P} \left[d_\infty(\hat{\mathbf{d}}_{k,n}(\mathbf{X}, \Theta), \mathcal{A}_k^*(\mathbf{X}, \Theta)) \leq \xi \right] \geq 1 - \rho.$$

We are now ready to prove Proposition 2. Fix $\rho, \xi > 0$. Since almost sure convergence implies convergence in probability, according to Lemma 1, there exists $k_0 \in \mathbb{N}^*$ such that

$$\mathbb{P} \left[\Delta(m, A_{k_0}^*(\mathbf{X}, \Theta)) \leq \xi \right] \geq 1 - \rho. \quad (3)$$

By Lemma 3, for all $\xi_1 > 0$, there exists $N \in \mathbb{N}^*$ such that, for all $n \geq N$,

$$\mathbb{P} \left[d_\infty(\hat{\mathbf{d}}_{k_0,n}(\mathbf{X}, \Theta), \mathcal{A}_{k_0}^*(\mathbf{X}, \Theta)) \leq \xi_1 \right] \geq 1 - \rho. \quad (4)$$

Since m is uniformly continuous, we can choose ξ_1 sufficiently small such that, for all $\mathbf{x} \in [0, 1]^p$, for all $\mathbf{d}_{k_0}, \mathbf{d}'_{k_0}$ satisfying $d_\infty(\mathbf{d}_{k_0}, \mathbf{d}'_{k_0}) \leq \xi_1$, we have

$$|\Delta(m, A(\mathbf{x}, \mathbf{d}_{k_0})) - \Delta(m, A(\mathbf{x}, \mathbf{d}'_{k_0}))| \leq \xi. \quad (5)$$

Thus, combining inequalities (4) and (5), we obtain

$$\mathbb{P} \left[|\Delta(m, A_{k_0,n}(\mathbf{X}, \Theta)) - \Delta(m, A_{k_0}^*(\mathbf{X}, \Theta))| \leq \xi \right] \geq 1 - \rho. \quad (6)$$

Using the fact that $\Delta(m, A) \leq \Delta(m, A')$ whenever $A \subset A'$, we deduce from (3) and (6) that, for all $n \geq N$,

$$\mathbb{P} [\Delta(m, A_n(\mathbf{X}, \Theta)) \leq 2\xi] \geq 1 - 2\rho.$$

This concludes the proof of Proposition 2.

5.3 Proof of Theorem 3.1

We still need some additional notations. The partition obtained with the random variable Θ and the data set \mathcal{D}_n is denoted by $\mathcal{P}_n(\mathcal{D}_n, \Theta)$, which we abbreviate as $\mathcal{P}_n(\Theta)$. We let

$$\Pi_n(\Theta) = \{\mathcal{P}((\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n), \Theta) : (\mathbf{x}_i, y_i) \in [0, 1]^d \times \mathbb{R}\}$$

be the family of all achievable partitions with random parameter Θ . Accordingly, we let

$$M(\Pi_n(\Theta)) = \max \{\text{Card}(\mathcal{P}) : \mathcal{P} \in \Pi_n(\Theta)\}$$

be the maximal number of terminal nodes among all partitions in $\Pi_n(\Theta)$. Given a set $\mathbf{z}_1^n = \{\mathbf{z}_1, \dots, \mathbf{z}_n\} \subset [0, 1]^d$, $\Gamma(\mathbf{z}_1^n, \Pi_n(\Theta))$ denotes the number of distinct partitions of \mathbf{z}_1^n induced by elements of $\Pi_n(\Theta)$, that is, the number of different partitions $\{\mathbf{z}_1^n \cap A : A \in \mathcal{P}\}$ of \mathbf{z}_1^n , for $\mathcal{P} \in \Pi_n(\Theta)$. Consequently, the partitioning number $\Gamma_n(\Pi_n(\Theta))$ is defined by

$$\Gamma_n(\Pi_n(\Theta)) = \max \{\Gamma(\mathbf{z}_1^n, \Pi_n(\Theta)) : \mathbf{z}_1, \dots, \mathbf{z}_n \in [0, 1]^d\}.$$

Let $(\beta_n)_n$ be a positive sequence, and define the truncated operator T_{β_n} by

$$\begin{cases} T_{\beta_n} u = u & \text{if } |u| < \beta_n \\ T_{\beta_n} u = \text{sign}(u)\beta_n & \text{if } |u| \geq \beta_n. \end{cases}$$

Hence, $T_{\beta_n} m_n(\mathbf{X}, \Theta)$, $Y_L = T_L Y$ and $Y_{i,L} = T_L Y_i$ are defined unambiguously. We let $\mathcal{F}(\Theta)$ be the set of all functions $f : [0, 1]^d \rightarrow \mathbb{R}$ piecewise constant on each cell of the partition $\mathcal{P}_n(\Theta)$. Finally, we denote by $\mathcal{I}_{n,\Theta}$ the set of indices of the data points that are selected during the subsampling step. Thus the tree estimate $m_n(\mathbf{x}, \Theta)$ satisfies

$$m_n(\cdot, \Theta) \in \underset{f \in \mathcal{F}_n(\Theta)}{\text{argmin}} \frac{1}{a_n} \sum_{i \in \mathcal{I}_{n,\Theta}} |f(\mathbf{X}_i) - Y_i|^2.$$

The proof of Theorem 3.1 is based on ideas developed by Nobel (1996), and worked out in Theorem 10.2 in Györfi et al. (2002). This theorem, tailored for our context, is recalled below for the sake of completeness.

Theorem 5.1. (*Györfi et al., 2002*) Let m_n and $\mathcal{F}_n(\Theta)$ be as above. Assume that

$$(i) \quad \lim_{n \rightarrow \infty} \beta_n = \infty,$$

$$(ii) \quad \lim_{n \rightarrow \infty} \mathbb{E} \left[\inf_{\substack{f \in \mathcal{F}_n(\Theta) \\ \|f\|_\infty \leq \beta_n}} \mathbb{E}_{\mathbf{X}} [f(\mathbf{X}) - m(\mathbf{X})]^2 \right] = 0,$$

(iii) For all $L > 0$,

$$\lim_{n \rightarrow \infty} \mathbb{E} \left[\sup_{\substack{f \in \mathcal{F}_n(\Theta) \\ \|f\|_\infty \leq \beta_n}} \left| \frac{1}{a_n} \sum_{i \in \mathcal{I}_{n,\Theta}} [f(\mathbf{X}_i) - Y_{i,L}]^2 - \mathbb{E}[f(\mathbf{X}) - Y_L]^2 \right| \right] = 0.$$

Then

$$\lim_{n \rightarrow \infty} \mathbb{E} [T_{\beta_n} m_n(\mathbf{X}, \Theta) - m(\mathbf{X})]^2 = 0.$$

We are now equipped to prove Theorem 3.1. Fix $\xi > 0$ and note that we just have to check statements (i) – (iii) of Theorem 5.1 to prove that the truncated estimate of the random forest is consistent. Throughout the proof, we let $\beta_n = \|m\|_\infty + \sigma\sqrt{2}(\log a_n)^2$. Clearly, statement (i) is true. To prove (ii), let

$$f_{n,\Theta} = \sum_{A \in \mathcal{P}_n(\Theta)} m(\mathbf{z}_A) \mathbb{1}_A,$$

where $\mathbf{z}_A \in A$ is an arbitrary point picked in cell A . Since, according to **(H1)**, $\|m\|_\infty < \infty$, for all n large enough such that $\beta_n > \|m\|_\infty$, we have

$$\begin{aligned} \mathbb{E} \inf_{\substack{f \in \mathcal{F}_n(\Theta) \\ \|f\|_\infty \leq \beta_n}} \mathbb{E}_{\mathbf{X}} [f(\mathbf{X}) - m(\mathbf{X})]^2 &\leq \mathbb{E} \inf_{\substack{f \in \mathcal{F}_n(\Theta) \\ \|f\|_\infty \leq \|m\|_\infty}} \mathbb{E}_{\mathbf{X}} [f(\mathbf{X}) - m(\mathbf{X})]^2 \\ &\leq \mathbb{E} [f_{\Theta,n}(\mathbf{X}) - m(\mathbf{X})]^2 \\ &\quad (\text{since } f_{\Theta,n} \in \mathcal{F}_n(\Theta)) \\ &\leq \mathbb{E} [m(\mathbf{z}_{A_n(\mathbf{X}, \Theta)}) - m(\mathbf{X})]^2 \\ &\leq \mathbb{E} [\Delta(m, A_n(\mathbf{X}, \Theta))]^2 \\ &\leq \xi^2 + 4\|m\|_\infty^2 \mathbb{P}[\Delta(m, A_n(\mathbf{X}, \Theta)) > \xi]. \end{aligned}$$

Thus, using Proposition 2, we see that, for all n large enough,

$$\mathbb{E} \inf_{\substack{f \in \mathcal{F}_n(\Theta) \\ \|f\|_\infty \leq \beta_n}} \mathbb{E}_{\mathbf{X}} [f(\mathbf{X}) - m(\mathbf{X})]^2 \leq 2\xi^2.$$

This establishes (ii).

To prove statement (iii), fix $L > 0$. Then, for all n large enough such that $L < \beta_n$,

$$\begin{aligned} & \mathbb{P}_{\mathbf{X}, \mathcal{D}_n} \left(\sup_{\substack{f \in \mathcal{F}_n(\Theta) \\ \|f\|_\infty \leq \beta_n}} \left| \frac{1}{a_n} \sum_{i \in \mathcal{I}_{n, \Theta}} [f(\mathbf{X}_i) - Y_{i,L}]^2 - \mathbb{E}[f(\mathbf{X}) - Y_L]^2 \right| > \xi \right) \\ & \leq 8 \exp \left[\log \Gamma_n(\Pi_n(\Theta)) + 2M(\Pi_n(\Theta)) \log \left(\frac{333e\beta_n^2}{\xi} \right) - \frac{a_n \xi^2}{2048\beta_n^4} \right] \\ & \quad (\text{according to Theorem 9.1 in Györfi et al., 2002}) \\ & \leq 8 \exp \left[-\frac{a_n}{\beta_n^4} \left(\frac{\xi^2}{2048} - \frac{\beta_n^4 \log \Gamma_n(\Pi_n)}{a_n} - \frac{2\beta_n^4 M(\Pi_n)}{a_n} \log \left(\frac{333e\beta_n^2}{\xi} \right) \right) \right]. \end{aligned}$$

Since each tree has exactly t_n terminal nodes, we have $M(\Pi_n(\Theta)) = t_n$ and simple calculations show that

$$\Gamma_n(\Pi_n(\Theta)) \leq (da_n)^{t_n}.$$

Hence,

$$\begin{aligned} & \mathbb{P} \left(\sup_{\substack{f \in \mathcal{F}_n(\Theta) \\ \|f\|_\infty \leq \beta_n}} \left| \frac{1}{a_n} \sum_{i \in \mathcal{I}_{n, \Theta}} [f(\mathbf{X}_i) - Y_{i,L}]^2 - \mathbb{E}[f(\mathbf{X}) - Y_L]^2 \right| > \xi \right) \\ & \leq 8 \exp \left(-\frac{a_n C_{\xi, n}}{\beta_n^4} \right), \end{aligned}$$

where

$$\begin{aligned} C_{\xi, n} &= \frac{\xi^2}{2048} - 4\sigma^4 \frac{t_n (\log(da_n))^9}{a_n} - 8\sigma^4 \frac{t_n (\log a_n)^8}{a_n} \log \left(\frac{666e\sigma^2 (\log a_n)^4}{\xi} \right) \\ & \rightarrow \frac{\xi^2}{2048}, \quad \text{as } n \rightarrow \infty, \end{aligned}$$

by our assumption. Finally, observe that

$$\sup_{\substack{f \in \mathcal{F}_n(\Theta) \\ \|f\|_\infty \leq \beta_n}} \left| \frac{1}{a_n} \sum_{i \in \mathcal{I}_{n, \Theta}} [f(\mathbf{X}_i) - Y_{i,L}]^2 - \mathbb{E}[f(\mathbf{X}) - Y_L]^2 \right| \leq 2(\beta_n + L)^2,$$

which yields, for all n large enough,

$$\begin{aligned}
& \mathbb{E} \left[\sup_{\substack{f \in \mathcal{F}_n(\Theta) \\ \|f\|_\infty \leq \beta_n}} \left| \frac{1}{a_n} \sum_{i=1}^{a_n} [f(\mathbf{X}_i) - Y_{i,L}]^2 - \mathbb{E}[f(\mathbf{X}) - Y_L]^2 \right| \right] \leq \xi \\
& \quad + 2(\beta_n + L)^2 \mathbb{P} \left[\sup_{\substack{f \in \mathcal{F}_n(\Theta) \\ \|f\|_\infty \leq \beta_n}} \left| \frac{1}{a_n} \sum_{i=1}^{a_n} [f(\mathbf{X}_i) - Y_{i,L}]^2 - \mathbb{E}[f(\mathbf{X}) - Y_L]^2 \right| > \xi \right] \\
& \leq \xi + 16(\beta_n + L)^2 \exp \left(-\frac{a_n C_{\xi,n}}{\beta_n^4} \right) \\
& \leq 2\xi.
\end{aligned}$$

Thus, according to Theorem 5.1,

$$\mathbb{E}[T_{\beta_n} m_n(\mathbf{X}, \Theta) - m(\mathbf{X})]^2 \rightarrow 0.$$

It remains to show the consistency of the non truncated random forest estimate, and the proof will be complete. For that purpose, note that, for all n large enough,

$$\begin{aligned}
\mathbb{E}[m_n(\mathbf{X}) - m(\mathbf{X})]^2 &= \mathbb{E}[\mathbb{E}_\Theta[m_n(\mathbf{X}, \Theta)] - m(\mathbf{X})]^2 \\
&\leq \mathbb{E}[m_n(\mathbf{X}, \Theta) - m(\mathbf{X})]^2 \\
&\quad (\text{by Jensen's inequality}) \\
&\leq \mathbb{E}[m_n(\mathbf{X}, \Theta) - T_{\beta_n} m_n(\mathbf{X}, \Theta)]^2 \\
&\quad + \mathbb{E}[T_{\beta_n} m_n(\mathbf{X}, \Theta) - m(\mathbf{X})]^2 \\
&\leq \mathbb{E} \left[[m_n(\mathbf{X}, \Theta) - T_{\beta_n} m_n(\mathbf{X}, \Theta)]^2 \mathbf{1}_{m_n(\mathbf{x}, \Theta) \geq \beta_n} \right] + \xi \\
&\leq \mathbb{E} \left[m_n^2(\mathbf{X}, \Theta) \mathbf{1}_{m_n(\mathbf{x}, \Theta) \geq \beta_n} \right] + \xi \\
&\leq \mathbb{E} \left[\mathbb{E} \left[m_n^2(\mathbf{X}, \Theta) \mathbf{1}_{m_n(\mathbf{x}, \Theta) \geq \beta_n} \mid \Theta \right] \right] + \xi.
\end{aligned}$$

Since $|m_n(\mathbf{X}, \Theta)| \leq \|m\|_\infty + \max_{1 \leq i \leq n} |\varepsilon_i|$, we have

$$\begin{aligned}
& \mathbb{E} \left[m_n^2(\mathbf{X}, \Theta) \mathbf{1}_{m_n(\mathbf{x}, \Theta) \geq \beta_n} \mid \Theta \right] \\
& \leq \mathbb{E} \left[(2\|m\|_\infty^2 + 2 \max_{1 \leq i \leq a_n} \varepsilon_i^2) \mathbf{1}_{\max_{1 \leq i \leq a_n} \varepsilon_i \geq \sigma\sqrt{2}(\log a_n)^2} \right] \\
& \leq 2\|m\|_\infty^2 \mathbb{P} \left[\max_{1 \leq i \leq a_n} \varepsilon_i \geq \sigma\sqrt{2}(\log a_n)^2 \right] \\
& \quad + 2 \left(\mathbb{E} \left[\max_{1 \leq i \leq a_n} \varepsilon_i^4 \right] \mathbb{P} \left[\max_{1 \leq i \leq a_n} \varepsilon_i \geq \sigma\sqrt{2}(\log a_n)^2 \right] \right)^{1/2}.
\end{aligned}$$

It is easy to see that

$$\mathbb{P}\left[\max_{1 \leq i \leq a_n} \varepsilon_i \geq \sigma\sqrt{2}(\log a_n)^2\right] \leq \frac{a_n^{1-\log a_n}}{2\sqrt{\pi}(\log a_n)^2}.$$

Finally, since the ε_i 's are centered i.i.d. Gaussian random variables, we have, for all n large enough,

$$\begin{aligned} \mathbb{E}[m_n(\mathbf{X}) - m(\mathbf{X})]^2 &\leq \frac{2\|m\|_\infty^2 a_n^{1-\log a_n}}{2\sqrt{\pi}(\log a_n)^2} + \xi + 2\left(3a_n\sigma^4 \frac{a_n^{1-\log a_n}}{2\sqrt{\pi}(\log a_n)^2}\right)^{1/2} \\ &\leq 3\xi. \end{aligned}$$

This concludes the proof of Theorem 3.1.

5.4 Proof of Theorem 3.2

Recall that each cell contains exactly one data point. Thus, letting

$$W_{ni}(\mathbf{X}) = \mathbb{E}_\Theta [\mathbf{1}_{\mathbf{X}_i \in A_n(\mathbf{X}, \Theta)}],$$

the random forest estimate m_n may be rewritten as

$$m_n(\mathbf{X}) = \sum_{i=1}^n W_{ni}(\mathbf{X}) Y_i.$$

We have in particular that $\sum_{i=1}^n W_{ni}(\mathbf{X}) = 1$. Thus,

$$\begin{aligned} \mathbb{E}[m_n(\mathbf{X}) - m(\mathbf{X})]^2 &\leq 2\mathbb{E}\left[\sum_{i=1}^n W_{ni}(\mathbf{X})(Y_i - m(\mathbf{X}_i))\right]^2 \\ &\quad + 2\mathbb{E}\left[\sum_{i=1}^n W_{ni}(\mathbf{X})(m(\mathbf{X}_i) - m(\mathbf{X}))\right]^2 \\ &\stackrel{\text{def}}{=} 2I_n + 2J_n. \end{aligned}$$

Fix $\alpha > 0$. To upper bound J_n , note that by Jensen's inequality,

$$\begin{aligned} J_n &\leq \mathbb{E}\left[\sum_{i=1}^n \mathbf{1}_{\mathbf{X}_i \in A_n(\mathbf{X}, \Theta)} (m(\mathbf{X}_i) - m(\mathbf{X}))^2\right] \\ &\leq \mathbb{E}\left[\sum_{i=1}^n \mathbf{1}_{\mathbf{X}_i \in A_n(\mathbf{X}, \Theta)} \Delta^2(m, A_n(\mathbf{X}, \Theta))\right] \\ &\leq \mathbb{E}[\Delta^2(m, A_n(\mathbf{X}, \Theta))]. \end{aligned}$$

So, by definition of $\Delta(m, A_n(\mathbf{X}, \Theta))^2$,

$$\begin{aligned} J_n &\leq 4\|m\|_\infty^2 \mathbb{E}[\mathbb{1}_{\Delta^2(m, A_n(\mathbf{X}, \Theta)) \geq \alpha}] + \alpha \\ &\leq \alpha(4\|m\|_\infty^2 + 1), \end{aligned}$$

for all n large enough, according to Proposition 2.

To bound I_n from above, we note that

$$\begin{aligned} I_n &= \mathbb{E} \left[\sum_{i,j=1}^n W_{ni}(\mathbf{X})W_{nj}(\mathbf{X})(Y_i - m(\mathbf{X}_i))(Y_j - m(\mathbf{X}_j)) \right] \\ &= \mathbb{E} \left[\sum_{i=1}^n W_{ni}^2(\mathbf{X})(Y_i - m(\mathbf{X}_i))^2 \right] + I'_n, \end{aligned}$$

where

$$I'_n = \mathbb{E} \left[\sum_{\substack{i,j \\ i \neq j}} \mathbb{1}_{\mathbf{x} \overset{\Theta}{\leftrightarrow} \mathbf{X}_i} \mathbb{1}_{\mathbf{x} \overset{\Theta'}{\leftrightarrow} \mathbf{X}_j} (Y_i - m(\mathbf{X}_i))(Y_j - m(\mathbf{X}_j)) \right].$$

The term I'_n , which involves the double products, is handled separately in Lemma 4 below. According to this lemma, and by assumption **(H2)**, for all n large enough,

$$|I'_n| \leq \alpha.$$

Consequently, recalling that $\varepsilon_i = Y_i - m(\mathbf{X}_i)$, we have, for all n large enough,

$$\begin{aligned} |I_n| &\leq \alpha + \mathbb{E} \left[\sum_{i=1}^n W_{ni}^2(\mathbf{X})(Y_i - m(\mathbf{X}_i))^2 \right] \\ &\leq \alpha + \mathbb{E} \left[\max_{1 \leq \ell \leq n} W_{n\ell}(\mathbf{X}) \sum_{i=1}^n W_{ni}(\mathbf{X}) \varepsilon_i^2 \right] \\ &\leq \alpha + \mathbb{E} \left[\max_{1 \leq \ell \leq n} W_{n\ell}(\mathbf{X}) \max_{1 \leq i \leq n} \varepsilon_i^2 \right]. \end{aligned} \tag{7}$$

Now, observe that in the subsampling step, there are exactly $\binom{a_n-1}{n-1}$ choices to pick a fixed observation \mathbf{X}_i . Since \mathbf{x} and \mathbf{X}_i belong to the same cell only if \mathbf{X}_i is selected in the subsampling step, we see that

$$\mathbb{P}_\Theta \left[\mathbf{X} \overset{\Theta}{\leftrightarrow} \mathbf{X}_i \right] \leq \frac{\binom{a_n-1}{n-1}}{\binom{a_n}{n}} = \frac{a_n}{n},$$

where \mathbb{P}_Θ denotes the probability with respect to Θ , conditional on \mathbf{X} and \mathcal{D}_n . So,

$$\max_{1 \leq i \leq n} W_{ni}(\mathbf{X}) \leq \max_{1 \leq i \leq n} \mathbb{P}_\Theta \left[\mathbf{X} \overset{\Theta}{\leftrightarrow} \mathbf{X}_i \right] \leq \frac{a_n}{n}. \quad (8)$$

Thus, combining inequalities (7) and (8), for all n large enough,

$$|I_n| \leq \alpha + \frac{a_n}{n} \mathbb{E} \left[\max_{1 \leq i \leq n} \varepsilon_i^2 \right].$$

The term inside the brackets is the maximum of n χ^2 -squared distributed random variables. Thus, for some positive constant C ,

$$\mathbb{E} \left[\max_{1 \leq i \leq n} \varepsilon_i^2 \right] \leq C \log n$$

(see, e.g., Chapter 1 in [Boucheron et al., 2013](#)). We conclude that, for all n large enough,

$$I_n \leq \alpha + C \frac{a_n \log n}{n} \leq 2\alpha.$$

Since α was arbitrary, the proof is complete.

Lemma 4. *Assume that **(H2)** is satisfied. Then, for all $\varepsilon > 0$, and all n large enough, $|I'_n| \leq \alpha$.*

Proof of Lemma 4. Firstly, assume that **(H2.2)** holds. Thus, we have for all $\ell_1, \ell_2 \in \{0, 1\}$,

$$\begin{aligned} & \text{Corr}(Y_i - m(\mathbf{X}_i), \mathbb{1}_{Z_{i,j}=(\ell_1, \ell_2)} | \mathbf{X}_i, \mathbf{X}_j, Y_j) \\ &= \frac{\mathbb{E}[(Y_i - m(\mathbf{X}_i)) \mathbb{1}_{Z_{i,j}=(\ell_1, \ell_2)}]}{\mathbb{V}^{1/2}[Y_i - m(\mathbf{X}_i) | \mathbf{X}_i, \mathbf{X}_j, Y_j] \mathbb{V}^{1/2}[\mathbb{1}_{Z_{i,j}=(\ell_1, \ell_2)} | \mathbf{X}_i, \mathbf{X}_j, Y_j]} \\ &= \frac{\mathbb{E}[(Y_i - m(\mathbf{X}_i)) \mathbb{1}_{Z_{i,j}=(\ell_1, \ell_2)} | \mathbf{X}_i, \mathbf{X}_j, Y_j]}{\sigma(\mathbb{P}[Z_{i,j} = (\ell_1, \ell_2) | \mathbf{X}_i, \mathbf{X}_j, Y_j] - \mathbb{P}[Z_{i,j} = (\ell_1, \ell_2) | \mathbf{X}_i, \mathbf{X}_j, Y_j]^2)^{1/2}} \\ &\geq \frac{\mathbb{E}[(Y_i - m(\mathbf{X}_i)) \mathbb{1}_{Z_{i,j}=(\ell_1, \ell_2)} | \mathbf{X}_i, \mathbf{X}_j, Y_j]}{\sigma \mathbb{P}^{1/2}[Z_{i,j} = (\ell_1, \ell_2) | \mathbf{X}_i, \mathbf{X}_j, Y_j]}, \end{aligned}$$

where the first equality comes from the fact that, for all $\ell_1, \ell_2 \in \{0, 1\}$,

$$\begin{aligned} & \text{Cov}(Y_i - m(\mathbf{X}_i), \mathbb{1}_{Z_{i,j}=(\ell_1, \ell_2)} | \mathbf{X}_i, \mathbf{X}_j, Y_j) \\ &= \mathbb{E}[(Y_i - m(\mathbf{X}_i)) \mathbb{1}_{Z_{i,j}=(\ell_1, \ell_2)} | \mathbf{X}_i, \mathbf{X}_j, Y_j], \end{aligned}$$

since $\mathbb{E}[Y_i - m(\mathbf{X}_i)|\mathbf{X}_i, \mathbf{X}_j, Y_j] = 0$. Thus, noticing that, almost surely,

$$\begin{aligned} & \mathbb{E} \left[Y_i - m(\mathbf{X}_i) \middle| Z_{i,j}, \mathbf{X}_i, \mathbf{X}_j, Y_j \right] \\ &= \sum_{\ell_1, \ell_2=1}^2 \frac{\mathbb{E} \left[(Y_i - m(\mathbf{X}_i)) \mathbf{1}_{Z_{i,j}=(\ell_1, \ell_2)} | \mathbf{X}_i, \mathbf{X}_j, Y_j \right]}{\mathbb{P} \left[Z_{i,j} = (\ell_1, \ell_2) | \mathbf{X}_i, \mathbf{X}_j, Y_j \right]} \mathbf{1}_{Z_{i,j}=(\ell_1, \ell_2)} \\ &\leq 4\sigma \max_{\ell_1, \ell_2=0,1} \frac{|\text{Corr}(Y_i - m(\mathbf{X}_i), \mathbf{1}_{Z_{i,j}=(\ell_1, \ell_2)} | \mathbf{X}_i, \mathbf{X}_j, Y_j)|}{\mathbb{P}^{1/2} \left[Z_{i,j} = (\ell_1, \ell_2) | \mathbf{X}_i, \mathbf{X}_j, Y_j \right]} \\ &\leq 4\sigma\gamma_n, \end{aligned}$$

we conclude that the first statement in **(H2.2)** implies that, almost surely,

$$\mathbb{E} \left[Y_i - m(\mathbf{X}_i) \middle| Z_{i,j}, \mathbf{X}_i, \mathbf{X}_j, Y_j \right] \leq 4\sigma\gamma_n.$$

Similarly, one can prove that second statement in assumption **(H2.2)** implies that, almost surely,

$$\mathbb{E} \left[|Y_i - m(\mathbf{X}_i)|^2 \middle| \mathbf{X}_i, \mathbf{1}_{\mathbf{X} \ni \mathbf{X}_i} \right] \leq 4C\sigma^2.$$

Returning to the term I'_n , and recalling that $W_{ni}(\mathbf{X}) = \mathbb{E}_{\Theta}[\mathbf{1}_{\mathbf{X} \ni \mathbf{X}_i}]$, we obtain

$$\begin{aligned} I'_n &= \mathbb{E} \left[\sum_{\substack{i,j \\ i \neq j}} \mathbf{1}_{\mathbf{X} \ni \mathbf{X}_i} \mathbf{1}_{\mathbf{X} \ni \mathbf{X}_j} (Y_i - m(\mathbf{X}_i))(Y_j - m(\mathbf{X}_j)) \right] \\ &= \sum_{\substack{i,j \\ i \neq j}} \mathbb{E} \left[\mathbb{E} \left[\mathbf{1}_{\mathbf{X} \ni \mathbf{X}_i} \mathbf{1}_{\mathbf{X} \ni \mathbf{X}_j} (Y_i - m(\mathbf{X}_i))(Y_j - m(\mathbf{X}_j)) \right. \right. \\ &\quad \left. \left. \middle| \mathbf{X}_i, \mathbf{X}_j, Y_i, \mathbf{1}_{\mathbf{X} \ni \mathbf{X}_i}, \mathbf{1}_{\mathbf{X} \ni \mathbf{X}_j} \right] \right] \\ &= \sum_{\substack{i,j \\ i \neq j}} \mathbb{E} \left[\mathbf{1}_{\mathbf{X} \ni \mathbf{X}_i} \mathbf{1}_{\mathbf{X} \ni \mathbf{X}_j} (Y_i - m(\mathbf{X}_i)) \right. \\ &\quad \left. \times \mathbb{E} \left[Y_j - m(\mathbf{X}_j) | \mathbf{X}_i, \mathbf{X}_j, Y_i, \mathbf{1}_{\mathbf{X} \ni \mathbf{X}_i}, \mathbf{1}_{\mathbf{X} \ni \mathbf{X}_j} \right] \right]. \end{aligned}$$

Therefore, by assumption **(H2.2)**,

$$\begin{aligned}
|I'_n| &\leq 4\sigma\gamma_n \sum_{\substack{i,j \\ i \neq j}} \mathbb{E} \left[\mathbf{1}_{\mathbf{X} \overset{\ominus}{\leftrightarrow} \mathbf{X}_i} \mathbf{1}_{\mathbf{X} \overset{\ominus'}{\leftrightarrow} \mathbf{X}_j} |Y_i - m(\mathbf{X}_i)| \right] \\
&\leq \gamma_n \sum_{i=1}^n \mathbb{E} \left[\mathbf{1}_{\mathbf{X} \overset{\ominus}{\leftrightarrow} \mathbf{X}_i} |Y_i - m(\mathbf{X}_i)| \right] \\
&\leq \gamma_n \sum_{i=1}^n \mathbb{E} \left[\mathbf{1}_{\mathbf{X} \overset{\ominus}{\leftrightarrow} \mathbf{X}_i} \mathbb{E} \left[|Y_i - m(\mathbf{X}_i)| \middle| \mathbf{X}_i, \mathbf{1}_{\mathbf{X} \overset{\ominus}{\leftrightarrow} \mathbf{X}_i} \right] \right] \\
&\leq \gamma_n \sum_{i=1}^n \mathbb{E} \left[\mathbf{1}_{\mathbf{X} \overset{\ominus}{\leftrightarrow} \mathbf{X}_i} \mathbb{E}^{1/2} \left[|Y_i - m(\mathbf{X}_i)|^2 \middle| \mathbf{X}_i, \mathbf{1}_{\mathbf{X} \overset{\ominus}{\leftrightarrow} \mathbf{X}_i} \right] \right] \\
&\leq 2\sigma C^{1/2} \gamma_n.
\end{aligned}$$

This proves the result, provided **(H2.2)** is true. Let us now assume that **(H2.1)** is verified. The key argument is to note that a data point \mathbf{X}_i can be connected with a random point \mathbf{X} if (\mathbf{X}_i, Y_i) is selected via the subsampling procedure and if there is no other data points in the hyperrectangle defined by \mathbf{X}_i and \mathbf{X} . Data points \mathbf{X}_i satisfying the latter geometrical property are called Layered Nearest Neighbor (LNN, see, e.g., [Barndorff-Nielsen and Sobel, 1966](#)). The connection between LNN and random forests has been first observed by [Lin and Jeon \(2006\)](#), and latter worked out by [Biau and Devroye \(2010\)](#). It is known, in particular, that the number of LNN $L_{a_n}(\mathbf{X})$ among a_n data points uniformly distributed on $[0, 1]^d$ satisfies, for some constant $C_1 > 0$ and for all n large enough,

$$\begin{aligned}
\mathbb{E}[L_{a_n}^4(\mathbf{X})] &\leq a_n \mathbb{P}[\mathbf{X} \overset{\ominus}{\underset{\text{LNN}}{\leftrightarrow}} \mathbf{X}_j] + 16a_n^2 \mathbb{P}[\mathbf{X} \overset{\ominus}{\underset{\text{LNN}}{\leftrightarrow}} \mathbf{X}_i] \mathbb{P}[\mathbf{X} \overset{\ominus}{\underset{\text{LNN}}{\leftrightarrow}} \mathbf{X}_j] \\
&\leq C_1 (\log a_n)^{2d-2},
\end{aligned} \tag{9}$$

(see, e.g., [Barndorff-Nielsen and Sobel, 1966](#); [Bai et al., 2005](#)). Thus, we have

$$I'_n = \mathbb{E} \left[\sum_{\substack{i,j \\ i \neq j}} \mathbf{1}_{\mathbf{X} \overset{\ominus}{\leftrightarrow} \mathbf{X}_i} \mathbf{1}_{\mathbf{X} \overset{\ominus'}{\leftrightarrow} \mathbf{X}_j} \mathbf{1}_{\mathbf{X}_i \overset{\ominus}{\underset{\text{LNN}}{\leftrightarrow}} \mathbf{X}} \mathbf{1}_{\mathbf{X}_j \overset{\ominus'}{\underset{\text{LNN}}{\leftrightarrow}} \mathbf{X}} (Y_i - m(\mathbf{X}_i))(Y_j - m(\mathbf{X}_j)) \right].$$

Consequently,

$$I'_n = \mathbb{E} \left[\sum_{\substack{i,j \\ i \neq j}} (Y_i - m(\mathbf{X}_i))(Y_j - m(\mathbf{X}_j)) \mathbb{1}_{\mathbf{X}_i \xrightarrow{\Theta} \mathbf{X}} \mathbb{1}_{\mathbf{X}_j \xrightarrow{\Theta'} \mathbf{X}} \right. \\ \left. \times \mathbb{E} \left[\mathbb{1}_{\mathbf{X} \xrightarrow{\Theta} \mathbf{X}_i} \mathbb{1}_{\mathbf{X} \xrightarrow{\Theta'} \mathbf{X}_j} \mid \mathbf{X}, \Theta, \Theta', \mathbf{X}_1, \dots, \mathbf{X}_n, Y_i, Y_j \right] \right],$$

where $\mathbf{X}_i \xrightarrow{\Theta} \mathbf{X}$ is the event where \mathbf{X}_i is selected by the subsampling and is also a LNN of \mathbf{X} . Next, with notations of assumption **(H2)**,

$$I'_n = \mathbb{E} \left[\sum_{\substack{i,j \\ i \neq j}} (Y_i - m(\mathbf{X}_i))(Y_j - m(\mathbf{X}_j)) \mathbb{1}_{\mathbf{X}_i \xrightarrow{\Theta} \mathbf{X}} \mathbb{1}_{\mathbf{X}_j \xrightarrow{\Theta'} \mathbf{X}} \right. \\ \left. \times \psi_{i,j}(Y_i, Y_j) \right] \\ = \mathbb{E} \left[\sum_{\substack{i,j \\ i \neq j}} (Y_i - m(\mathbf{X}_i))(Y_j - m(\mathbf{X}_j)) \mathbb{1}_{\mathbf{X}_i \xrightarrow{\Theta} \mathbf{X}} \mathbb{1}_{\mathbf{X}_j \xrightarrow{\Theta'} \mathbf{X}} \psi_{i,j} \right] \\ + \mathbb{E} \left[\sum_{\substack{i,j \\ i \neq j}} (Y_i - m(\mathbf{X}_i))(Y_j - m(\mathbf{X}_j)) \mathbb{1}_{\mathbf{X}_i \xrightarrow{\Theta} \mathbf{X}} \mathbb{1}_{\mathbf{X}_j \xrightarrow{\Theta'} \mathbf{X}} \right. \\ \left. \times (\psi_{i,j}(Y_i, Y_j) - \psi_{i,j}) \right].$$

The first term is easily seen to be zero since

$$\mathbb{E} \left[\sum_{\substack{i,j \\ i \neq j}} (Y_i - m(\mathbf{X}_i))(Y_j - m(\mathbf{X}_j)) \mathbb{1}_{\mathbf{X}_i \xrightarrow{\Theta} \mathbf{X}} \mathbb{1}_{\mathbf{X}_j \xrightarrow{\Theta'} \mathbf{X}} \psi(\mathbf{X}, \Theta, \Theta', \mathbf{X}_1, \dots, \mathbf{X}_n) \right] \\ = \sum_{\substack{i,j \\ i \neq j}} \mathbb{E} \left[\mathbb{1}_{\mathbf{X}_i \xrightarrow{\Theta} \mathbf{X}} \mathbb{1}_{\mathbf{X}_j \xrightarrow{\Theta'} \mathbf{X}} \psi_{i,j} \right. \\ \left. \times \mathbb{E}[(Y_i - m(\mathbf{X}_i))(Y_j - m(\mathbf{X}_j)) \mid \mathbf{X}, \mathbf{X}_1, \dots, \mathbf{X}_n, \Theta, \Theta'] \right] \\ = 0.$$

Therefore,

$$\begin{aligned}
|I'_n| &\leq \mathbb{E} \left[\sum_{\substack{i,j \\ i \neq j}} |Y_i - m(\mathbf{X}_i)| |Y_j - m(\mathbf{X}_j)| \mathbb{1}_{\mathbf{X}_i \overset{\Theta}{\leftarrow} \mathbf{X}} \mathbb{1}_{\mathbf{X}_j \overset{\Theta'}{\leftarrow} \mathbf{X}} \right. \\
&\quad \left. \times |\psi_{i,j}(Y_i, Y_j) - \psi_{i,j}| \right] \\
&\leq \mathbb{E} \left[\max_{1 \leq \ell \leq n} |Y_\ell - m(\mathbf{X}_\ell)|^2 \max_{\substack{i,j \\ i \neq j}} |\psi_{i,j}(Y_i, Y_j) - \psi_{i,j}| \right. \\
&\quad \left. \times \sum_{\substack{i,j \\ i \neq j}} \mathbb{1}_{\mathbf{X}_i \overset{\Theta}{\leftarrow} \mathbf{X}} \mathbb{1}_{\mathbf{X}_j \overset{\Theta'}{\leftarrow} \mathbf{X}} \right].
\end{aligned}$$

Now, observe that

$$\sum_{\substack{i,j \\ i \neq j}} \mathbb{1}_{\mathbf{X}_i \overset{\Theta}{\leftarrow} \mathbf{X}} \mathbb{1}_{\mathbf{X}_j \overset{\Theta'}{\leftarrow} \mathbf{X}} \leq L_{a_n}^2(\mathbf{X}),$$

Consequently,

$$\begin{aligned}
|I'_n| &\leq \mathbb{E}^{1/2} \left[L_{a_n}^4(\mathbf{X}) \max_{1 \leq \ell \leq n} |Y_\ell - m(\mathbf{X}_\ell)|^4 \right] \\
&\quad \times \mathbb{E}^{1/2} \left[\max_{\substack{i,j \\ i \neq j}} |\psi_{i,j}(Y_i, Y_j) - \psi_{i,j}| \right]^2. \tag{10}
\end{aligned}$$

Simple calculations reveal that there exists $C_1 > 0$ such that, for all n ,

$$\mathbb{E} \left[\max_{1 \leq \ell \leq n} |Y_\ell - m(\mathbf{X}_\ell)|^4 \right] \leq C_1 (\log n)^2. \tag{11}$$

Thus, by inequalities (9) and (11), the first term in (10) can be upper bounded as follows:

$$\begin{aligned}
&\mathbb{E}^{1/2} \left[L_{a_n}^4(\mathbf{X}) \max_{1 \leq \ell \leq n} |Y_\ell - m(\mathbf{X}_\ell)|^4 \right] \\
&= \mathbb{E}^{1/2} \left[L_{a_n}^4(\mathbf{X}) \mathbb{E} \left[\max_{1 \leq \ell \leq n} |Y_\ell - m(\mathbf{X}_\ell)|^4 \mid \mathbf{X}, \mathbf{X}_1, \dots, \mathbf{X}_n \right] \right] \\
&\leq C' (\log n) (\log a_n)^{d-1}.
\end{aligned}$$

Finally,

$$|I'_n| \leq C'(\log a_n)^{d-1}(\log n)^{\alpha/2}\mathbb{E}^{1/2} \left[\max_{\substack{i,j \\ i \neq j}} |\psi_{i,j}(Y_i, Y_j) - \psi_{i,j}| \right]^2,$$

which tends to zero by assumption. □

6 Technical results

6.1 Proof of Lemma 1

Technical Lemma 1. *Assume that (H1) is satisfied and that $L^* \equiv 0$ for all cuts in some given cell A . Then the regression function m is constant on A .*

Proof of Technical Lemma 1. We start by proving the result in dimension $p = 1$. Letting $A = [a, b]$ ($0 \leq a < b \leq 1$), and recalling that $Y = m(\mathbf{X}) + \varepsilon$, one has

$$\begin{aligned} L^*(1, z) &= \mathbb{V}[Y|\mathbf{X} \in A] - \mathbb{P}[a \leq \mathbf{X} \leq z | \mathbf{X} \in A] \mathbb{V}[Y|a \leq \mathbf{X} \leq z] \\ &\quad - \mathbb{P}[z \leq \mathbf{X} \leq b | \mathbf{X} \in A] \mathbb{V}[Y|z < \mathbf{X} \leq b] \\ &= -\frac{1}{(b-a)^2} \left(\int_a^b m(t) dt \right)^2 + \frac{1}{(b-a)(z-a)} \left(\int_a^z m(t) dt \right)^2 \\ &\quad + \frac{1}{(b-a)(b-z)} \left(\int_z^b m(t) dt \right)^2. \end{aligned}$$

Let $C = \int_a^b m(t) dt$ and $M(z) = \int_a^z m(t) dt$. Simple calculations show that

$$L^*(1, z) = \frac{1}{(z-a)(b-z)} \left(M(z) - C \frac{z-a}{b-a} \right)^2.$$

Therefore, since $L^* \equiv 0$ on \mathcal{C}_A by assumption, we obtain

$$M(z) = C \frac{z-a}{b-a}.$$

This proves that $M(z)$ is linear in z , and that m is therefore constant on $[a, b]$.

Let us now examine the general multivariate case, where $A = \prod_{j=1}^p [a_j, b_j] \subset [0, 1]^p$. From the univariate analysis, we know that, for all $1 \leq j \leq p$, there exists a constant C_j such that

$$\int_{a_1}^{b_1} \dots \int_{a_p}^{b_p} m(\mathbf{x}) dx_1 \dots dx_{j-1} dx_{j+1} \dots dx_p = C_j.$$

Since m is additive this implies that, for all j and all x_j ,

$$m_j(x_j) = C_j - \int_{a_1}^{b_1} \dots \int_{a_p}^{b_p} \sum_{\ell \neq j} m_\ell(x_\ell) dx_1 \dots dx_{j-1} dx_{j+1} \dots dx_p,$$

which does not depend upon x_i . This shows that m is constant on A . \square

Proof of Lemma 1. Take $\xi > 0$ and fix $\mathbf{x} \in [0, 1]^p$. Let θ be a realization of the random variable Θ . Since m is uniformly continuous, the result is clear if $\text{diam}(A_k^*(\mathbf{x}, \theta))$ tends to zero as k tends to infinity. Thus, in the sequel, it is assumed that $\text{diam}(A_k^*(\mathbf{x}, \theta))$ does not tend to zero. In that case, since $(A_k^*(\mathbf{x}, \theta))_k$ is a decreasing sequence of compact sets, there exist $\mathbf{a}_\infty(\mathbf{x}, \theta) = (\mathbf{a}_\infty^{(1)}(\mathbf{x}, \theta), \dots, \mathbf{a}_\infty^{(p)}(\mathbf{x}, \theta)) \in [0, 1]^p$ and $\mathbf{b}_\infty(\mathbf{x}, \theta) = (\mathbf{b}_\infty^{(1)}(\mathbf{x}, \theta), \dots, \mathbf{b}_\infty^{(p)}(\mathbf{x}, \theta)) \in [0, 1]^p$ such that

$$\begin{aligned} \bigcap_{k=1}^{\infty} A_k^*(\mathbf{x}, \theta) &= \prod_{j=1}^p [\mathbf{a}_\infty^{(j)}(\mathbf{x}, \theta), \mathbf{b}_\infty^{(j)}(\mathbf{x}, \theta)] \\ &\stackrel{\text{def}}{=} A_\infty^*(\mathbf{x}, \theta). \end{aligned}$$

Since $\text{diam}(A_k^*(\mathbf{x}, \theta))$ does not tend to zero, there exists an index j' such that $\mathbf{a}_\infty^{(j')}(\mathbf{x}, \theta) < \mathbf{b}_\infty^{(j')}(\mathbf{x}, \theta)$ (i.e., the cell $A_\infty^*(\mathbf{x}, \theta)$ is not reduced to one point). Let $A_k^*(\mathbf{x}, \theta) \stackrel{\text{def}}{=} \prod_{j=1}^p [\mathbf{a}_k^{(j)}(\mathbf{x}, \theta), \mathbf{b}_k^{(j)}(\mathbf{x}, \theta)]$ be the cell containing \mathbf{x} at level k . If the criterion L^* is identically zero for all cuts in $A_\infty^*(\mathbf{x}, \theta)$ then m is constant on $A_\infty^*(\mathbf{x}, \theta)$ according to Lemma 1. This implies that $\Delta(m, A_\infty^*(\mathbf{x}, \theta)) = 0$. Thus, in that case, since m is uniformly continuous,

$$\lim_{k \rightarrow \infty} \Delta(m, A_k^*(\mathbf{x}, \theta)) = \Delta(m, A_\infty^*(\mathbf{x}, \theta)) = 0.$$

Let us now show by contradiction that L^* is almost surely necessarily null on the cuts of $A_\infty^*(\mathbf{x}, \theta)$. In the rest of the proof, for all $k \in \mathbb{N}^*$, we let L_k^* be

the criterion L^* used in the cell $A_k^*(\mathbf{x}, \theta)$, that is

$$\begin{aligned} L_k^*(d) &= \mathbb{V}[Y|\mathbf{X} \in A_k^*(\mathbf{x}, \theta)] \\ &\quad - \mathbb{P}[\mathbf{X}^{(j)} < z | \mathbf{X} \in A_k^*(\mathbf{x}, \theta)] \mathbb{V}[Y|\mathbf{X}^{(j)} < z, \mathbf{X} \in A_k^*(\mathbf{x}, \theta)] \\ &\quad - \mathbb{P}[\mathbf{X}^{(j)} \geq z | \mathbf{X} \in A_k^*(\mathbf{x}, \theta)] \mathbb{V}[Y|\mathbf{X}^{(j)} \geq z, \mathbf{X} \in A_k^*(\mathbf{x}, \theta)], \end{aligned}$$

for all $d = (j, z) \in \mathcal{C}_{A_k^*(\mathbf{x}, \theta)}$. If L_∞^* is not identically zero, then there exists a cut $d_\infty(\mathbf{x}, \theta)$ in $\mathcal{C}_{A_\infty^*(\mathbf{x}, \theta)}$ such that $L^*(d_\infty(\mathbf{x}, \theta)) = c > 0$. Fix $\xi > 0$. By the uniform continuity of m , there exists $\delta_1 > 0$ such that

$$\sup_{\|\mathbf{w} - \mathbf{w}'\|_\infty \leq \delta_1} |m(\mathbf{w}) - m(\mathbf{w}')| \leq \xi.$$

Since $A_k^*(\mathbf{x}, \theta) \downarrow A_\infty^*(\mathbf{x}, \theta)$, there exists k_0 such that, for all $k \geq k_0$,

$$\max(\|\mathbf{a}_k(\mathbf{x}, \theta) - \mathbf{a}_\infty(\mathbf{x}, \theta)\|_\infty, \|\mathbf{b}_k(\mathbf{x}, \theta) - \mathbf{b}_\infty(\mathbf{x}, \theta)\|_\infty) \leq \delta_1. \quad (12)$$

Observe that, for all $k \in \mathbb{N}^*$, $\mathbb{V}[Y|\mathbf{X} \in A_{k+1}^*(\mathbf{x}, \theta)] < \mathbb{V}[Y|\mathbf{X} \in A_k^*(\mathbf{x}, \theta)]$. Thus,

$$\underline{L}_k^* := \sup_{\substack{d \in \mathcal{C}_{A_k^*(\mathbf{x}, \theta)} \\ d^{(1)} \in \mathcal{M}_{\text{try}}} L_k^*(d) \leq \xi. \quad (13)$$

From inequality (12), we deduce that

$$|\mathbb{E}[m(\mathbf{X})|\mathbf{X} \in A_k^*(\mathbf{x}, \theta)] - \mathbb{E}[m(\mathbf{X})|\mathbf{X} \in A_\infty^*(\mathbf{x}, \theta)]| \leq \xi.$$

Consequently, there exists a constant $C > 0$ such that, for all $k \geq k_0$ and all cuts $d \in \mathcal{C}_{A_\infty^*(\mathbf{x}, \theta)}$,

$$|L_k^*(d) - L_\infty^*(d)| \leq C\xi^2. \quad (14)$$

Let $k_1 \geq k_0$ be the first level after k_0 at which the direction $d_\infty^{(1)}(\mathbf{x}, \theta)$ is amongst the m_{try} selected coordinates. Almost surely, $k_1 < \infty$. Thus, by the definition of $d_\infty(\mathbf{x}, \theta)$ and inequality (14),

$$c - C\xi^2 \leq L_\infty^*(d_\infty(\mathbf{x}, \theta)) - C\xi^2 \leq L_k^*(d_\infty(\mathbf{x}, \theta)),$$

which implies that $c - C\xi^2 \leq \underline{L}_k^*$. Hence, using inequality (13), we have

$$c - C\xi^2 \leq \underline{L}_k^* \leq \xi,$$

which is absurd, since $c > 0$ is fixed and ξ is arbitrarily small. Thus, by Lemma 1, m is constant on $A_\infty^*(\mathbf{x}, \theta)$. This implies that $\Delta(m, A_k^*(\mathbf{x}, \theta)) \rightarrow 0$ as $k \rightarrow \infty$.

6.2 Proof of Lemma 2

We start by proving Lemma 2 in the case $k = 1$, i.e., when the first cut is performed at the root of a tree. Since in that case $L_{n,1}(\mathbf{x}, \cdot)$ does not depend on \mathbf{x} , we simply write $L_{n,1}(\cdot)$ instead of $L_{n,1}(\mathbf{x}, \cdot)$.

Proof of Lemma 2 in the case $k = 1$. Fix $\alpha, \rho > 0$. Observe that if two cuts d_1, d_2 satisfy $\|d_1 - d_2\|_\infty < 1$, then the cut directions are the same, i.e., $d_1^{(1)} = d_2^{(1)}$. Using this fact and symmetry arguments, we just need to prove Lemma 2 when the cuts are performed along the first dimension. In other words, we only need to prove that

$$\lim_{n \rightarrow \infty} \mathbb{P} \left[\sup_{|x_1 - x_2| \leq \delta} |L_{n,1}(1, x_1) - L_{n,1}(1, x_2)| > \alpha \right] \leq \rho/p. \quad (15)$$

Letting $Z_i = \max_{1 \leq i \leq n} |\varepsilon_i|$, simple calculations show that

$$\mathbb{P}[Z_i \geq t] = 1 - \exp \left(n \ln (1 - 2\mathbb{P}[\varepsilon_1 \geq t]) \right).$$

The last probability can be upper bounded by using the following standard inequality on Gaussian tail:

$$\mathbb{P}[\varepsilon_1 \geq t] \leq \frac{\sigma}{t\sqrt{2\pi}} \exp \left(-\frac{t^2}{2\sigma^2} \right).$$

Consequently, there exists a constant $C_\rho > 0$ and $N_1 \in \mathbb{N}^*$ such that, with probability $1 - \rho$, for all $n > N_1$,

$$\max_{1 \leq i \leq n} |\varepsilon_i| \leq C_\rho \sqrt{\log n}. \quad (16)$$

Besides, by simple calculations on Gaussian tail, for all $n \in \mathbb{N}^*$, we have

$$\mathbb{P} \left[\left| \frac{1}{n} \sum_{i=1}^n \varepsilon_i \right| \geq \alpha \right] \leq \frac{\sigma}{\alpha\sqrt{n}} \exp \left(-\frac{\alpha^2 n}{2\sigma^2} \right).$$

Since there are, at most, n^2 sets of the form $\{i : X_i \in [a_n, b_n]\}$ for $0 \leq a_n < b_n \leq 1$, we deduce from the last inequality and the union bound, that there exists $N_2 \in \mathbb{N}^*$ such that, with probability $1 - \rho$, for all $n > N_2$ and all $0 \leq a_n < b_n \leq 1$ satisfying $N_n([a_n, b_n] \times [0, 1]^{p-1}) > \sqrt{n}$,

$$\left| \frac{1}{N_n([a_n, b_n] \times [0, 1]^{p-1})} \sum_{\substack{i: X_i \in [a_n, b_n] \\ \times [0, 1]^{p-1}}} \varepsilon_i \right| \leq \alpha. \quad (17)$$

By the Glivenko-Cantelli theorem, there exists $N_3 \in \mathbb{N}^*$ such that, with probability $1 - \rho$, for all $0 \leq a < b \leq 1$, and all $n > N_3$,

$$(b - a - \delta^2)n \leq N_n([a, b] \times [0, 1]^{p-1}) \leq (b - a + \delta^2)n. \quad (18)$$

Throughout the proof, we assume to be on the event where assertions (16)-(18) hold, which occurs with probability $1 - 3\rho$, for all $n > N$, where $N = \max(N_1, N_2, N_3)$.

Take $x_1, x_2 \in [0, 1]$ such that $|x_1 - x_2| \leq \delta$ and assume, without loss of generality, that $x_1 < x_2$. In the remainder of the proof, we will need the following quantities (see Figure 1 for an illustration in dimension two):

$$\begin{cases} A_{L, \sqrt{\delta}} = [0, \sqrt{\delta}] \times [0, 1]^{p-1} \\ A_{R, \sqrt{\delta}} = [1 - \sqrt{\delta}, 1] \times [0, 1]^{p-1} \\ A_{C, \sqrt{\delta}} = [\sqrt{\delta}, 1 - \sqrt{\delta}] \times [0, 1]^{p-1}. \end{cases}$$

Similarly, we define

$$\begin{cases} A_{L,1} = [0, x_1] \times [0, 1]^{p-1} \\ A_{R,1} = [x_1, 1] \times [0, 1]^{p-1} \\ A_{L,2} = [0, x_2] \times [0, 1]^{p-1} \\ A_{R,2} = [x_2, 1] \times [0, 1]^{p-1} \\ A_C = [x_1, x_2] \times [0, 1]^{p-1}. \end{cases}$$

Recall that, for any cell A , \bar{Y}_A is the mean of the Y_i 's falling in A and $N_n(A)$ is the number of data points in A . To prove (15), five cases are to be considered, depending upon the positions of x_1 and x_2 . We repeatedly use the decomposition

$$L_{n,1}(1, x_1) - L_{n,1}(1, x_2) = J_1 + J_2 + J_3,$$

where

$$\begin{aligned} J_1 &= \frac{1}{n} \sum_{i: \mathbf{X}_i^{(1)} < x_1} (Y_i - \bar{Y}_{A_{L,1}})^2 - \frac{1}{n} \sum_{i: \mathbf{X}_i^{(1)} < x_1} (Y_i - \bar{Y}_{A_{L,2}})^2, \\ J_2 &= \frac{1}{n} \sum_{i: \mathbf{X}_i^{(1)} \in [x_1, x_2]} (Y_i - \bar{Y}_{A_{R,1}})^2 - \frac{1}{n} \sum_{i: \mathbf{X}_i^{(1)} \in [x_1, x_2]} (Y_i - \bar{Y}_{A_{L,2}})^2, \\ \text{and } J_3 &= \frac{1}{n} \sum_{i: \mathbf{X}_i^{(1)} \geq x_2} (Y_i - \bar{Y}_{A_{R,1}})^2 - \frac{1}{n} \sum_{i: \mathbf{X}_i^{(1)} \geq x_2} (Y_i - \bar{Y}_{A_{R,2}})^2. \end{aligned}$$

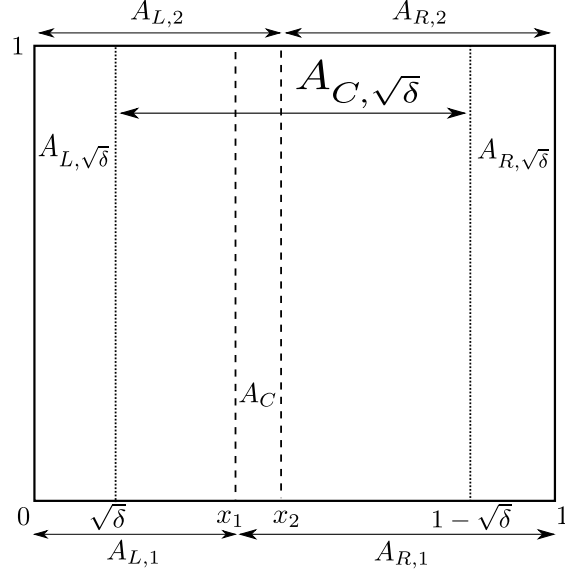


Figure 1: Illustration of the notation in dimension $p = 2$.

First case: $x_1, x_2 \in A_{C, \sqrt{\delta}}$. Since $N_n(A_{L,2}) > N_n(A_{L, \sqrt{\delta}}) > \sqrt{n}$ for all $n > N$, we have, according to inequalities (17),

$$|\bar{Y}_{A_{L,2}}| \leq \|m\|_\infty + \alpha \quad \text{and} \quad |\bar{Y}_{A_{R,1}}| \leq \|m\|_\infty + \alpha.$$

Therefore

$$\begin{aligned} |J_2| &= 2 |\bar{Y}_{A_{L,2}} - \bar{Y}_{A_{R,1}}| \times \frac{1}{n} \left| \sum_{i: \mathbf{X}_i^{(1)} \in [x_1, x_2]} \left(Y_i - \frac{\bar{Y}_{A_{L,2}} + \bar{Y}_{A_{R,1}}}{2} \right) \right| \\ &\leq 4(\|m\|_\infty + \alpha) \left(\frac{(\|m\|_\infty + \alpha)N_n(A_C)}{n} + \frac{1}{n} \left| \sum_{i: \mathbf{X}_i^{(1)} \in [x_1, x_2]} m(\mathbf{X}_i) \right| \right. \\ &\quad \left. + \frac{1}{n} \left| \sum_{i: \mathbf{X}_i^{(1)} \in [x_1, x_2]} \varepsilon_i \right| \right) \\ &\leq 4(\|m\|_\infty + \alpha) \left((\delta + \delta^2)(\|m\|_\infty + \alpha) + \|m\|_\infty(\delta + \delta^2) \right. \\ &\quad \left. + \frac{1}{n} \left| \sum_{i: \mathbf{X}_i^{(1)} \in [x_1, x_2]} \varepsilon_i \right| \right). \end{aligned}$$

If $N_n(A_C) \geq \sqrt{n}$, we obtain

$$\frac{1}{n} \left| \sum_{i: \mathbf{X}_i^{(1)} \in [x_1, x_2]} \varepsilon_i \right| \leq \frac{1}{N_n(A_C)} \left| \sum_{i: \mathbf{X}_i^{(1)} \in [x_1, x_2]} \varepsilon_i \right| \leq \alpha \quad (\text{according to (17)})$$

or, if $N_n(A_C) < \sqrt{n}$, we have

$$\frac{1}{n} \left| \sum_{i: \mathbf{X}_i^{(1)} \in [x_1, x_2]} \varepsilon_i \right| \leq \frac{C_\rho \sqrt{\log n}}{\sqrt{n}} \quad (\text{according to (16)}).$$

Thus, for all n large enough,

$$|J_2| \leq 4(\|m\|_\infty + \alpha) \left((\delta + \delta^2)(2\|m\|_\infty + \alpha) + \alpha \right). \quad (19)$$

With respect to J_1 , observe that

$$\begin{aligned} |\bar{Y}_{A_{L,1}} - \bar{Y}_{A_{L,2}}| &= \left| \frac{1}{N_n(A_{L,1})} \sum_{i: \mathbf{X}_i^{(1)} < x_1} Y_i - \frac{1}{N_n(A_{L,2})} \sum_{i: \mathbf{X}_i^{(1)} < x_2} Y_i \right| \\ &\leq \left| \frac{1}{N_n(A_{L,1})} \sum_{i: \mathbf{X}_i^{(1)} < x_1} Y_i - \frac{1}{N_n(A_{L,2})} \sum_{i: \mathbf{X}_i^{(1)} < x_1} Y_i \right| \\ &\quad + \left| \frac{1}{N_n(A_{L,2})} \sum_{i: \mathbf{X}_i^{(1)} \in [x_1, x_2]} Y_i \right| \\ &\leq \left| 1 - \frac{N_n(A_{L,1})}{N_n(A_{L,2})} \right| \times \frac{1}{N_n(A_{L,1})} \times \left| \sum_{i: \mathbf{X}_i^{(1)} < x_1} Y_i \right| \\ &\quad + \frac{1}{N_n(A_{L,2})} \left| \sum_{i: \mathbf{X}_i^{(1)} \in [x_1, x_2]} Y_i \right|. \end{aligned}$$

Since $N_n(A_{L,2}) - N_n(A_{L,1}) \leq n(\delta + \delta^2)$, we obtain

$$1 - \frac{N_n(A_{L,1})}{N_n(A_{L,2})} \leq \frac{n(\delta + \delta^2)}{N_n(A_{L,2})} \leq \frac{\delta + \delta^2}{\sqrt{\delta} - \delta^2} \leq 4\sqrt{\delta},$$

for all δ small enough, which implies that

$$\begin{aligned}
|\bar{Y}_{A_{L,1}} - \bar{Y}_{A_{L,2}}| &\leq \frac{4\sqrt{\delta}}{N_n(A_{L,1})} \left| \sum_{i:\mathbf{X}_i^{(1)} < x_1} Y_i \right| \\
&\quad + \frac{N_n(A_{L,1})}{N_n(A_{L,2})} \times \frac{1}{N_n(A_{L,1})} \left| \sum_{i:\mathbf{X}_i^{(1)} \in [x_1, x_2]} Y_i \right| \\
&\leq 4\sqrt{\delta}(\|m\|_\infty + \alpha) + \frac{N_n(A_{L,1})}{N_n(A_{L,2})}(\|m\|_\infty \delta + \alpha) \\
&\leq 5(\|m\|_\infty \sqrt{\delta} + \alpha).
\end{aligned}$$

Thus,

$$\begin{aligned}
|J_1| &= \left| \frac{1}{n} \sum_{i:\mathbf{X}_i^{(1)} < x_1} (Y_i - \bar{Y}_{A_{L,1}})^2 - \frac{1}{n} \sum_{i:\mathbf{X}_i^{(1)} < x_1} (Y_i - \bar{Y}_{A_{L,2}})^2 \right| \\
&= \left| (\bar{Y}_{A_{L,2}} - \bar{Y}_{A_{L,1}}) \times \frac{2}{n} \sum_{i:\mathbf{X}_i^{(1)} < x_1} \left(Y_i - \frac{\bar{Y}_{A_{L,1}} + \bar{Y}_{A_{L,2}}}{2} \right) \right| \\
&\leq |\bar{Y}_{A_{L,2}} - \bar{Y}_{A_{L,1}}|^2 \\
&\leq 25(\|m\|_\infty \sqrt{\delta} + \alpha)^2.
\end{aligned} \tag{20}$$

The term J_3 can be bounded with the same arguments.

Finally, by (19) and (20), for all $n > N$, and all δ small enough, we conclude that

$$\begin{aligned}
|L_n(1, x_1) - L_n(1, x_2)| &\leq 4(\|m\|_\infty + \alpha) \left((\delta + \delta^2)(2\|m\|_\infty + \alpha) + \alpha \right) \\
&\quad + 25(\|m\|_\infty \sqrt{\delta} + \alpha)^2 \\
&\leq \alpha.
\end{aligned}$$

Second case: $x_1, x_2 \in A_{L, \sqrt{\delta}}$. With the same arguments as above, one proves that

$$\begin{aligned}
|J_1| &\leq \max \left(4(\sqrt{\delta} + \delta^2)(\|m\|_\infty + \alpha)^2, \alpha \right), \\
|J_2| &\leq \max(4(\|m\|_\infty + \alpha)(2\delta\|m\|_\infty + 2\alpha), \alpha), \\
|J_3| &\leq 25(\|m\|_\infty \sqrt{\delta} + \alpha)^2.
\end{aligned}$$

Consequently, for all n large enough,

$$|L_n(1, x_1) - L_n(1, x_2)| = J_1 + J_2 + J_3 \leq 3\alpha.$$

The other cases $\{x_1, x_2 \in A_{R, \sqrt{\delta}}\}$, $\{x_1, x_2 \in A_{L, \sqrt{\delta}} \times A_{C, \sqrt{\delta}}\}$, and $\{x_1, x_2 \in A_{C, \sqrt{\delta}} \times A_{R, \sqrt{\delta}}\}$ can be treated in the same way. Details are omitted. \square

Proof of Lemma 2. We proceed similarly as in the proof of the case $k = 1$. Here, we establish the result for $k = 2$ and $p = 2$ only. Extensions are easy and left to the reader. Fix $\rho > 0$. At first, it should be noted that there exists $N_1 \in \mathbb{N}^*$ such that, with probability $1 - \rho$, for all $n > N_0$ and all $A_n = [a_n^{(1)}, b_n^{(1)}] \times [a_n^{(2)}, b_n^{(2)}] \subset [0, 1]^2$ satisfying $N_n(A_n) > \sqrt{n}$, we have

$$\left| \frac{1}{N_n(A_n)} \sum_{i: X_i \in A_n} \varepsilon_i \right| \leq \alpha, \quad (21)$$

and

$$\frac{1}{N_n(A_n)} \sum_{i: X_i \in A_n} \varepsilon_i^2 \leq \tilde{\sigma}^2, \quad (22)$$

where $\tilde{\sigma}^2$ is a positive constant, depending only on ρ . Inequality (22) is a straightforward consequence of the following inequality (see, e.g., [Laurent and Massart, 2000](#)), which is valid for all $n \in \mathbb{N}^*$:

$$\mathbb{P} [\chi^2(n) \geq 5n] \leq \exp(-n).$$

Throughout the proof, we assume to be on the event where assertions (16), (18), (21)-(22) hold, which occurs with probability $1 - 3\rho$, for all n large enough. We also assume that $d_1 = (1, x_1)$ and $d_2 = (2, x_2)$ (see Figure 2). The other cases can be treated similarly.

Let $d'_1 = (1, x'_1)$ and $d'_2 = (2, x'_2)$ be such that $|x_1 - x'_1| < \delta$ and $|x_2 - x'_2| < \delta$. Then the CART-split criterion $L_{n,2}$ writes

$$\begin{aligned} L_n(d_1, d_2) &= \frac{1}{N_n(A_{R,1})} \sum_i (Y_i - \bar{Y}_{A_{R,1}})^2 \mathbb{1}_{\mathbf{X}_i^{(1)} > x_1} \\ &\quad - \frac{1}{N_n(A_{R,1})} \sum_{i: \mathbf{X}_i^{(2)} > x_2} (Y_i - \bar{Y}_{A_{H,2}})^2 \mathbb{1}_{\mathbf{X}_i^{(1)} > x_1} \\ &\quad - \frac{1}{N_n(A_{R,1})} \sum_{i: \mathbf{X}_i^{(2)} \leq x_2} (Y_i - \bar{Y}_{A_{B,2}})^2 \mathbb{1}_{\mathbf{X}_i^{(1)} > x_1}. \end{aligned}$$

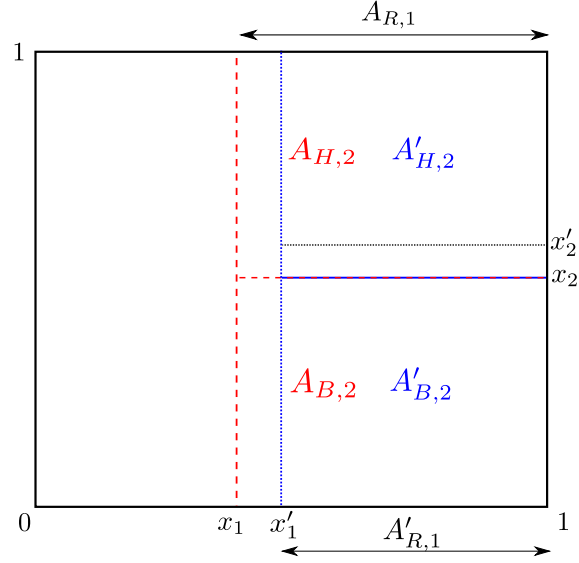


Figure 2: An example of cells in dimension $p = 2$.

Clearly,

$$L_n(d_1, d_2) - L_n(d'_1, d'_2) = L_n(d_1, d_2) - L_n(d'_1, d_2) + L_n(d'_1, d_2) - L_n(d'_1, d'_2).$$

We have (Figure 2):

$$\begin{aligned}
L_n(d_1, d_2) - L_n(d'_1, d_2) &= \left[\frac{1}{N_n(A_{R,1})} \sum_{i: \mathbf{X}_i^{(2)} > x_2} (Y_i - \bar{Y}_{A_{H,2}})^2 \mathbf{1}_{\mathbf{X}_i^{(1)} > x_1} \right. \\
&\quad \left. - \frac{1}{N_n(A'_{R,1})} \sum_{i: \mathbf{X}_i^{(2)} > x_2} (Y_i - \bar{Y}_{A'_{H,2}})^2 \mathbf{1}_{\mathbf{X}_i^{(1)} > x'_1} \right] \\
&+ \left[\frac{1}{N_n(A_{R,1})} \sum_{i: \mathbf{X}_i^{(2)} \leq x_2} (Y_i - \bar{Y}_{A_{B,2}})^2 \mathbf{1}_{\mathbf{X}_i^{(1)} > x_1} \right. \\
&\quad \left. - \frac{1}{N_n(A'_{R,1})} \sum_{i: \mathbf{X}_i^{(2)} \leq x_2} (Y_i - \bar{Y}_{A'_{B,2}})^2 \mathbf{1}_{\mathbf{X}_i^{(1)} > x'_1} \right] \\
&\stackrel{\text{def}}{=} A_1 + B_1.
\end{aligned}$$

The term A_1 can be rewritten as $A_1 = A_{1,1} + A_{1,2} + A_{1,3}$, where

$$\begin{aligned}
A_{1,1} &= \frac{1}{N_n(A_{R,1})} \sum_{i:\mathbf{X}_i^{(2)} > x_2} (Y_i - \bar{Y}_{A_{H,2}})^2 \mathbb{1}_{\mathbf{X}_i^{(1)} > x'_1} \\
&\quad - \frac{1}{N_n(A_{R,1})} \sum_{i:\mathbf{X}_i^{(2)} > x_2} (Y_i - \bar{Y}_{A'_{H,2}})^2 \mathbb{1}_{\mathbf{X}_i^{(1)} > x'_1}, \\
A_{1,2} &= \frac{1}{N_n(A_{R,1})} \sum_{i:\mathbf{X}_i^{(2)} > x_2} (Y_i - \bar{Y}_{A'_{H,2}})^2 \mathbb{1}_{\mathbf{X}_i^{(1)} > x'_1} \\
&\quad - \frac{1}{N_n(A'_{R,1})} \sum_{i:\mathbf{X}_i^{(2)} > x_2} (Y_i - \bar{Y}_{A'_{H,2}})^2 \mathbb{1}_{\mathbf{X}_i^{(1)} > x'_1}, \\
\text{and } A_{1,3} &= \frac{1}{N_n(A_{R,1})} \sum_{i:\mathbf{X}_i^{(2)} > x_2} (Y_i - \bar{Y}_{A_{H,2}})^2 \mathbb{1}_{\mathbf{X}_i^{(1)} \in [x_1, x'_1]}.
\end{aligned}$$

Easy calculations show that

$$A_{1,1} = \frac{N_n(A'_{H,2})}{N_n(A_{R,1})} (\bar{Y}_{A'_{H,2}} - \bar{Y}_{A_{H,2}})^2,$$

which implies, with the same arguments as in the proof for $k = 1$, that $A_{1,1} \rightarrow 0$ as $n \rightarrow \infty$. With respect to $A_{1,2}$ and $A_{1,3}$, we write

$$\max(A_{1,2}, A_{1,3}) \leq \max\left(C_\rho \frac{\log n}{\sqrt{n}}, 2(\tilde{\sigma}^2 + 4\|m\|_\infty^2 + \alpha^2) \frac{\sqrt{\delta}}{\xi}\right).$$

Thus, $A_{1,2} \rightarrow 0$ and $A_{1,3} \rightarrow 0$ as $n \rightarrow \infty$. Collecting bounds, we conclude that $A_1 \rightarrow 0$. One proves with similar arguments that $B_1 \rightarrow 0$ and, consequently, that $L_n(d'_1, d_2) - L_n(d'_1, d'_2) \rightarrow 0$. \square

6.3 Proof of Lemma 3

We prove by induction that, for all k , with probability $1 - \rho$, for all $\xi > 0$ and all n large enough,

$$d_\infty(\hat{\mathbf{d}}_{k,n}(\mathbf{X}, \Theta), \mathcal{A}_k^*(\mathbf{X}, \Theta)) \leq \xi.$$

Call this property H_k . Fix $k > 1$ and assume that H_{k-1} is true. For all $\mathbf{d}_{k-1} \in \mathcal{A}_{k-1}(\mathbf{X})$, let

$$\hat{d}_{k,n}(\mathbf{d}_{k-1}) \in \arg \min_{d_k} L_n(\mathbf{X}, \mathbf{d}_{k-1}, d_k),$$

and

$$d_k^*(\mathbf{d}_{k-1}) \in \arg \min_{d_k} L^*(\mathbf{X}, \mathbf{d}_{k-1}, d_k),$$

where the minimum is evaluated, as usual, over $\{d_k \in \mathcal{C}_{A(\mathbf{x}, \mathbf{d}_{k-1})} : d_k^{(1)} \in \mathcal{M}_{\text{try}}\}$. Fix $\rho > 0$. In the rest of the proof, we assume Θ to be fixed. Moreover, since Θ is fixed, we omit the dependence on Θ . We momentarily consider $\mathbf{x} \in [0, 1]^d$.

Note that, for all \mathbf{d}_{k-1} ,

$$\begin{aligned} & L_n(\mathbf{x}, \mathbf{d}_{k-1}, \hat{d}_{k,n}(\mathbf{d}_{k-1})) - L^*(\mathbf{x}, \mathbf{d}_{k-1}, \hat{d}_{k,n}(\mathbf{d}_{k-1})) \\ & \leq L_n(\mathbf{x}, \mathbf{d}_{k-1}, \hat{d}_{k,n}(\mathbf{d}_{k-1})) - L^*(\mathbf{x}, \mathbf{d}_{k-1}, d_k^*(\mathbf{d}_{k-1})) \\ & \quad (\text{by definition of } d_k^*(\mathbf{d}_{k-1})) \\ & \leq L_n(\mathbf{x}, \mathbf{d}_{k-1}, d_k^*(\mathbf{d}_{k-1})) - L^*(\mathbf{x}, \mathbf{d}_{k-1}, d_k^*(\mathbf{d}_{k-1})) \\ & \quad (\text{by definition of } \hat{d}_{k,n}(\mathbf{d}_{k-1})). \end{aligned}$$

Thus,

$$\begin{aligned} & \left| L_n(\mathbf{x}, \mathbf{d}_{k-1}, \hat{d}_{k,n}(\mathbf{d}_{k-1})) - L^*(\mathbf{x}, \mathbf{d}_{k-1}, d_k^*(\mathbf{d}_{k-1})) \right| \\ & \leq \max \left(\left| L_n(\mathbf{x}, \mathbf{d}_{k-1}, \hat{d}_{k,n}(\mathbf{d}_{k-1})) - L^*(\mathbf{x}, \mathbf{d}_{k-1}, \hat{d}_{k,n}(\mathbf{d}_{k-1})) \right|, \right. \\ & \quad \left. \left| L_n(\mathbf{x}, \mathbf{d}_{k-1}, d_k^*(\mathbf{d}_{k-1})) - L^*(\mathbf{x}, \mathbf{d}_{k-1}, d_k^*(\mathbf{d}_{k-1})) \right| \right) \\ & \leq \sup_{d_k} |L_n(\mathbf{x}, \mathbf{d}_{k-1}, d_k) - L^*(\mathbf{x}, \mathbf{d}_{k-1}, d_k)|. \end{aligned}$$

Moreover,

$$\begin{aligned} & |L^*(\mathbf{x}, \mathbf{d}_{k-1}, \hat{d}_{k,n}(\mathbf{d}_{k-1})) - L^*(\mathbf{x}, \mathbf{d}_{k-1}, d_k^*(\mathbf{d}_{k-1}))| \\ & \leq |L^*(\mathbf{x}, \mathbf{d}_{k-1}, \hat{d}_{k,n}(\mathbf{d}_{k-1})) - L_n(\mathbf{x}, \mathbf{d}_{k-1}, \hat{d}_{k,n}(\mathbf{d}_{k-1}))| \\ & \quad + |L_n(\mathbf{x}, \mathbf{d}_{k-1}, \hat{d}_{k,n}(\mathbf{d}_{k-1})) - L^*(\mathbf{x}, \mathbf{d}_{k-1}, d_k^*(\mathbf{d}_{k-1}))| \\ & \leq 2 \sup_{d_k} |L_n(\mathbf{x}, \mathbf{d}_{k-1}, d_k) - L^*(\mathbf{x}, \mathbf{d}_{k-1}, d_k)| \\ & = 2 \sup_{d_k} |L_n(\mathbf{x}, \mathbf{d}_k) - L^*(\mathbf{x}, \mathbf{d}_k)|. \end{aligned} \tag{23}$$

Let $\bar{\mathcal{A}}_k^\xi(\mathbf{x}) = \{\mathbf{d}_k : \mathbf{d}_{k-1} \in \mathcal{A}_{k-1}^\xi(\mathbf{x})\}$. So, taking the supremum on both sides of (23) leads to

$$\begin{aligned} & \sup_{\mathbf{d}_{k-1} \in \mathcal{A}_{k-1}^\xi(\mathbf{x})} |L^*(\mathbf{x}, \mathbf{d}_{k-1}, \hat{d}_{k,n}(\mathbf{d}_{k-1})) - L^*(\mathbf{x}, \mathbf{d}_{k-1}, d_k^*(\mathbf{d}_{k-1}))| \\ & \leq 2 \sup_{\mathbf{d}_k \in \bar{\mathcal{A}}_k^\xi(\mathbf{x})} |L_n(\mathbf{x}, \mathbf{d}_k) - L^*(\mathbf{x}, \mathbf{d}_k)|. \end{aligned} \tag{24}$$

By Lemma 2, for all $\xi' > 0$, one can find $\delta > 0$ such that, for all n large enough,

$$\mathbb{P} \left[\sup_{\mathbf{x} \in [0,1]^d} \sup_{\substack{\|\mathbf{d}_k - \mathbf{d}'_k\|_\infty \leq \delta \\ \mathbf{d}_k, \mathbf{d}'_k \in \bar{\mathcal{A}}_k^\xi(\mathbf{x})}} |L_n(\mathbf{x}, \mathbf{d}_k) - L_n(\mathbf{x}, \mathbf{d}'_k)| \leq \xi' \right] \geq 1 - \rho. \quad (25)$$

Now, let \mathcal{G} be a regular grid of $[0,1]^d$ whose grid step equal to $\xi/2$. Note that, for all $\mathbf{x} \in \mathcal{G}$, $\bar{\mathcal{A}}_k^\xi(\mathbf{x})$ is compact. Thus, for all $\mathbf{x} \in \mathcal{G}$, there exists a finite subset $\mathcal{C}_{\delta, \mathbf{x}} = \{c_{j, \mathbf{x}} : 1 \leq j \leq p\}$ such that, for all $\mathbf{d}_k \in \bar{\mathcal{A}}_k^\xi(\mathbf{x})$, $d_\infty(\mathbf{d}_k, \mathcal{C}_{\delta, \mathbf{x}}) \leq \delta$. Set $\xi' > 0$. Observe that, since the subset $\cup_{\mathbf{x} \in \mathcal{G}} \mathcal{C}_{\delta, \mathbf{x}}$ is finite, one has, for all n large enough,

$$\sup_{\mathbf{x} \in \mathcal{G}} \sup_{c_{j, \mathbf{x}} \in \mathcal{C}_{\delta, \mathbf{x}}} |L_n(\mathbf{x}, c_{j, \mathbf{x}}) - L^*(\mathbf{x}, c_{j, \mathbf{x}})| \leq \xi'. \quad (26)$$

Hence, for all n large enough,

$$\begin{aligned} \sup_{\mathbf{x} \in \mathcal{G}} \sup_{\mathbf{d}_k \in \bar{\mathcal{A}}_k^\xi(\mathbf{x})} |L_n(\mathbf{x}, \mathbf{d}_k) - L^*(\mathbf{x}, \mathbf{d}_k)| &\leq \sup_{\mathbf{x} \in \mathcal{G}} \sup_{\mathbf{d}_k \in \bar{\mathcal{A}}_k^\xi(\mathbf{x})} \left(|L_n(\mathbf{x}, \mathbf{d}_k) - L_n(\mathbf{x}, c_{j, \mathbf{x}})| \right. \\ &\quad \left. + |L_n(\mathbf{x}, c_{j, \mathbf{x}}) - L^*(\mathbf{x}, c_{j, \mathbf{x}})| + |L^*(\mathbf{x}, c_{j, \mathbf{x}}) - L^*(\mathbf{x}, \mathbf{d}_k)| \right), \end{aligned}$$

where $c_{j, \mathbf{x}}$ satisfies $\|c_{j, \mathbf{x}} - \mathbf{d}_k\|_\infty \leq \delta$. Using inequalities (25) and (26), with probability $1 - \rho$, we obtain, for all n large enough,

$$\sup_{\mathbf{x} \in \mathcal{G}} \sup_{\mathbf{d}_k \in \bar{\mathcal{A}}_k^\xi(\mathbf{x})} |L_n(\mathbf{x}, \mathbf{d}_k) - L^*(\mathbf{x}, \mathbf{d}_k)| \leq 3\xi'.$$

Finally, by inequality (24), with probability $1 - \rho$, for all n large enough,

$$\sup_{\mathbf{x} \in \mathcal{G}} \sup_{\mathbf{d}_{k-1} \in \mathcal{A}_{k-1}^\xi(\mathbf{x})} |L^*(\mathbf{x}, \mathbf{d}_{k-1}, \hat{d}_{k, n}(\mathbf{d}_{k-1})) - L^*(\mathbf{x}, \mathbf{d}_{k-1}, d_k^*(\mathbf{d}_{k-1}))| \leq 6\xi'. \quad (27)$$

Hereafter, to simplify, we assume that, for any given $(k-1)$ -tuple of theoretical cuts, there is only one theoretical cut at level k , and leave the general case as an easy adaptation. Thus, we can define unambiguously

$$d_k^*(\mathbf{d}_{k-1}) = \arg \min_{d_k} L^*(\mathbf{d}_{k-1}, d_k).$$

Fix $\xi'' > 0$. From inequality (27), by evoking the equicontinuity of L_n and the compactness of $\mathcal{U} = \{(\mathbf{x}, \mathbf{d}_{k-1}) : \mathbf{x} \in \mathcal{G}, \mathbf{d}_{k-1} \in \mathcal{A}_{k-1}^\xi(\mathbf{x})\}$, we deduce that, with probability $1 - \rho$, for all n large enough,

$$\sup_{(\mathbf{x}, \mathbf{d}_{k-1}) \in \mathcal{U}} d_\infty(\hat{d}_{k, n}(\mathbf{d}_{k-1}), d_k^*(\mathbf{d}_{k-1})) \leq \xi''. \quad (28)$$

Besides,

$$\mathbb{P}[(\mathbf{X}, \hat{\mathbf{d}}_{k-1,n}(\mathbf{X})) \in \mathcal{U}] = \mathbb{E}[\mathbb{P}[(\mathbf{X}, \hat{\mathbf{d}}_{k-1,n}(\mathbf{X})) \in \mathcal{U} | \mathcal{D}_n]] \geq 1 - 2^{k-1}\xi. \quad (29)$$

In the rest of the proof, we consider $\xi \leq \rho/2^{k-1}$, which, by inequalities (28) and (29), leads to

$$\mathbb{P}\left[\sup_{(\mathbf{x}, \mathbf{d}_{k-1}) \in \mathcal{U}} d_\infty(\hat{d}_{k,n}(\mathbf{d}_{k-1}), d_k^*(\mathbf{d}_{k-1})) \leq \xi'', (\mathbf{X}, \hat{\mathbf{d}}_{k-1,n}(\mathbf{X})) \in \mathcal{U}\right] \geq 1 - 2\rho.$$

This implies, with probability $1 - 2\rho$, for all n large enough,

$$d_\infty(\hat{d}_{k,n}(\hat{\mathbf{d}}_{k-1,n}), d_k^*(\hat{\mathbf{d}}_{k-1,n})) \leq \xi''. \quad (30)$$

Now, using triangle inequality,

$$\begin{aligned} d_\infty(\hat{d}_{k,n}(\hat{\mathbf{d}}_{k-1,n}), \mathcal{A}_k^*) &\leq d_\infty(\hat{d}_{k,n}(\hat{\mathbf{d}}_{k-1,n}), d_k^*(\hat{\mathbf{d}}_{k-1,n})) \\ &\quad + d_\infty(d_k^*(\hat{\mathbf{d}}_{k-1,n}), \mathcal{A}_k^*). \end{aligned} \quad (31)$$

Thus, we just have to show that $d_\infty(d_k^*(\hat{\mathbf{d}}_{k-1,n}), \mathcal{A}_k^*) \rightarrow 0$ in probability as $n \rightarrow \infty$, and the proof will be complete. To avoid confusion, we let $\{\mathbf{d}_{k-1}^{*,i} : i \in \mathcal{I}\}$ be the set of best first $(k-1)$ -th theoretical cuts (which can be either countable or not). With this notation, $d_k^*(\mathbf{d}_{k-1}^{*,i})$ is the k -th theoretical cuts given that the $(k-1)$ previous ones are $\mathbf{d}_{k-1}^{*,i}$. For simplicity, let

$$L^{i,*}(\mathbf{x}, d_k) = L_k^*(\mathbf{x}, \mathbf{d}_{k-1}^{*,i}, d_k) \quad \text{and} \quad \hat{L}^*(\mathbf{x}, d_k) = L_k^*(\mathbf{x}, \hat{\mathbf{d}}_{k-1,n}, d_k).$$

As before,

$$d_k^*(\mathbf{d}_{k-1}^{*,i}) \in \arg \min_{d_k} L^{i,*}(\mathbf{x}, d_k) \quad \text{and} \quad d_k^*(\hat{\mathbf{d}}_{k-1,n}) \in \arg \min_{d_k} \hat{L}^*(\mathbf{x}, d_k).$$

Clearly, the result will be proved if we establish that,

$$\inf_{i \in \mathcal{I}} d_\infty(d_k^*(\hat{\mathbf{d}}_{k-1,n}), d_k^*(\mathbf{d}_{k-1}^{*,i})) \rightarrow 0, \quad \text{in probability, as } n \rightarrow \infty.$$

Note that, for all $\mathbf{x} \in \mathcal{G}$, $\bar{\mathcal{A}}_k^{\xi}(\mathbf{x})$ is compact. Thus, for all $\mathbf{x} \in \mathcal{G}$, there exists a finite subset $\mathcal{C}'_{\delta, \mathbf{x}} = \{c'_{j, \mathbf{x}} : 1 \leq j \leq p\}$ such that, for all d_k , $d_\infty(d_k, \mathcal{C}'_{\delta, \mathbf{x}}) \leq \delta$. Hence, with probability $1 - \rho$, for all n large enough,

$$\begin{aligned} |\hat{L}^*(\mathbf{x}, d_k) - L^{i,*}(\mathbf{x}, d_k)| &\leq |\hat{L}^*(\mathbf{x}, d_k) - \hat{L}^*(\mathbf{x}, c'_{j, \mathbf{x}})| \\ &\quad + |\hat{L}^*(\mathbf{x}, c'_{j, \mathbf{x}}) - L^{i,*}(\mathbf{x}, c'_{j, \mathbf{x}})| \\ &\quad + |L^{i,*}(\mathbf{x}, c'_{j, \mathbf{x}}) - L^{i,*}(\mathbf{x}, d_k)| \\ &\leq 2\xi' + |\hat{L}^*(\mathbf{x}, c'_{j, \mathbf{x}}) - L^{i,*}(\mathbf{x}, c'_{j, \mathbf{x}})| \\ &\quad \text{(by the continuity of } L_k^* \text{)}. \end{aligned}$$

Therefore, as in inequality (24), with probability $1 - \rho$, for all i and all n large enough,

$$\begin{aligned} |L^{i,*}(\mathbf{x}, d_k^*(\hat{\mathbf{d}}_{k-1,n})) - L^{i,*}(\mathbf{x}, d_k^*(\mathbf{d}_{k-1}^{*,i}))| &\leq 2 \sup_{d_k} |\hat{L}^*(\mathbf{x}, d_k) - L^{i,*}(\mathbf{x}, d_k)| \\ &\leq 4\xi' + 2 \max_j |\hat{L}^*(\mathbf{x}, c'_{j,\mathbf{x}}) - L^{i,*}(\mathbf{x}, c'_{j,\mathbf{x}})|. \end{aligned}$$

Taking the infimum over all i , we obtain

$$\begin{aligned} \inf_i |L^{i,*}(\mathbf{x}, d_k^*(\hat{\mathbf{d}}_{k-1,n})) - L^{i,*}(\mathbf{x}, d_k^*(\mathbf{d}_{k-1}^{*,i}))| &\leq 4\xi' \\ &+ 2 \inf_i \max_j |\hat{L}^*(\mathbf{x}, c'_{j,\mathbf{x}}) - L^{i,*}(\mathbf{x}, c'_{j,\mathbf{x}})|. \end{aligned} \quad (32)$$

Introduce ω , the modulus of continuity of L_k^* :

$$\omega(\mathbf{x}, \delta) = \sup_{\|\mathbf{d} - \mathbf{d}'\|_\infty \leq \delta} |L_k^*(\mathbf{x}, \mathbf{d}) - L_k^*(\mathbf{x}, \mathbf{d}')|.$$

Observe that, since $L_k^*(\mathbf{x}, \cdot)$ is uniformly continuous, $\omega(\delta) \rightarrow 0$ as $\delta \rightarrow 0$. Hence, for all n large enough,

$$\begin{aligned} &\inf_i \max_j |\hat{L}^*(\mathbf{x}, c'_{j,\mathbf{x}}) - L^{i,*}(\mathbf{x}, c'_{j,\mathbf{x}})| \\ &= \inf_i \max_j |L_k^*(\mathbf{x}, \hat{\mathbf{d}}_{k-1,n}, c'_{j,\mathbf{x}}) - L_k^*(\mathbf{x}, \mathbf{d}_{k-1}^{*,i}, c'_{j,\mathbf{x}})| \\ &\leq \inf_i \omega(\mathbf{x}, \|\hat{\mathbf{d}}_{k-1,n} - \mathbf{d}_{k-1}^{*,i}\|_\infty) \\ &\leq \xi', \end{aligned} \quad (33)$$

since, by assumption H_{k-1} , $\inf_i \|\hat{\mathbf{d}}_{k-1,n} - \mathbf{d}_{k-1}^{*,i}\|_\infty \rightarrow 0$. Therefore, combining (32) and (33), with probability $1 - \rho$, for all n large enough,

$$\inf_i |L^{i,*}(\mathbf{X}, d_k^*(\hat{\mathbf{d}}_{k-1,n})) - L^{i,*}(\mathbf{X}, d_k^*(\mathbf{d}_{k-1}^{*,i}))| \leq 6\xi.$$

Finally, by Lemma 5 below, with probability $1 - \rho$, for all n large enough,

$$\inf_i d_\infty(d_k^*(\hat{\mathbf{d}}_{k-1,n}), d_k^*(\mathbf{d}_{k-1}^{*,i})) \leq \xi''. \quad (34)$$

Plugging inequality (34) and (30) into (31), we conclude that, with probability $1 - 3\rho$, for all n large enough,

$$d_\infty(\hat{d}_{k,n}(\hat{\mathbf{d}}_{k-1,n}), \mathcal{A}_k^*) \leq 2\xi'',$$

which proves H_k . Property H_1 can be proved in the same way.

Lemma 5. For all $\delta, \rho > 0$, there exists $\xi > 0$ such that, if, with probability $1 - \rho$,

$$\inf_i |L^{i,*}(\mathbf{X}, d_k^*(\hat{\mathbf{d}}_{k-1,n})) - L^{i,*}(\mathbf{X}, d_k^*(\mathbf{d}_{k-1}^{*,i}))| \leq \xi,$$

then, with probability $1 - \rho$,

$$\inf_i d_\infty(d_k^*(\hat{\mathbf{d}}_{k-1,n}), d_k^*(\mathbf{d}_{k-1}^{*,i})) \leq \delta. \quad (35)$$

Proof of Lemma 5. Fix $\rho > 0$. Note that, for all $\delta > 0$, there exists $\xi > 0$ such that,

$$\inf_{\mathbf{x} \in [0,1]^d} \inf_i \inf_{y: d_\infty(y, d_k^*(\mathbf{d}_{k-1}^{*,i})) \geq \delta} |L_k^*(\mathbf{x}, \mathbf{d}_{k-1}^{*,i}, d_k^*(\mathbf{d}_{k-1}^{*,i})) - L_k^*(\mathbf{x}, \mathbf{d}_{k-1}^{*,i}, y)| \geq \xi.$$

To see this, assume that one can find $\delta > 0$ such that, for all $\xi > 0$, there exist $i_\xi, y_\xi, \mathbf{x}_\xi$ satisfying

$$|L_k^*(\mathbf{x}_\xi, \mathbf{d}_{k-1}^{*,i_\xi}, d_k^*(\mathbf{d}_{k-1}^{*,i_\xi})) - L_k^*(\mathbf{x}_\xi, \mathbf{d}_{k-1}^{*,i_\xi}, y_\xi)| \leq \xi,$$

with $d_\infty(y_\xi, d_k^*(\mathbf{d}_{k-1}^{*,i_\xi})) \geq \delta$. Recall that $\{\mathbf{d}_{k-1}^{*,i} : i \in \mathbb{N}\}, \{d_k^*(\mathbf{d}_{k-1}^{*,i}) : i \in \mathbb{N}\}$ are compact. Then, letting $\xi_p = 1/p$, we can extract three sequences $\mathbf{d}_{k-1}^{*,i_p} \rightarrow \mathbf{d}_{k-1}, d_k^*(\mathbf{d}_{k-1}^{*,i_p}) \rightarrow d_k$ and $y_{\xi_{i_p}} \rightarrow y$ as $p \rightarrow \infty$ such that

$$L_k^*(\mathbf{d}_{k-1}, d_k) = L_k^*(\mathbf{d}_{k-1}, y), \quad (36)$$

and $d_\infty(y, d_k) \geq \delta$. Since we assume that given the $(k-1)$ -th first cuts \mathbf{d}_{k-1} , there is only one best cut d_k , equation (36) implies that $y = d_k$, which is absurd.

Now, to conclude the proof, fix $\delta > 0$ and assume that, with probability $1 - \rho$,

$$\inf_i d_\infty(d_k^*(\mathbf{d}_{k-1}^{*,i}), d_k^*(\hat{\mathbf{d}}_{k-1,n})) \geq \delta.$$

Thus, with probability $1 - \rho$,

$$\begin{aligned} & \inf_i |L^{i,*}(\mathbf{X}, d_k^*(\hat{\mathbf{d}}_{k-1,n})) - L^{i,*}(\mathbf{X}, d_k^*(\mathbf{d}_{k-1}^{*,i}))| \\ &= \inf_i |L_k^*(\mathbf{X}, \mathbf{d}_{k-1}^{*,i}, d_k^*(\hat{\mathbf{d}}_{k-1,n})) - L_k^*(\mathbf{X}, \mathbf{d}_{k-1}^{*,i}, d_k^*(\mathbf{d}_{k-1}^{*,i}))| \\ &\geq \inf_{\mathbf{x} \in [0,1]^d} \inf_i \inf_{d_\infty(y, d_k^*(\mathbf{d}_{k-1}^{*,i})) \geq \delta} |L_k^*(\mathbf{x}, \mathbf{d}_{k-1}^{*,i}, y) - L_k^*(\mathbf{x}, \mathbf{d}_{k-1}^{*,i}, d_k^*(\mathbf{d}_{k-1}^{*,i}))| \\ &\geq \xi, \end{aligned}$$

which, by contraposition, concludes the proof. \square

Proof of Proposition 1. Fix $k \in \mathbb{N}^*$ and $\rho, \xi > 0$. According to Lemma 3, with probability $1 - \rho$, for all n large enough, there exists a sequence of theoretical first k cuts $\mathbf{d}_k^*(\mathbf{X}, \Theta)$ such that

$$d_\infty(\mathbf{d}_k^*(\mathbf{X}, \Theta), \hat{\mathbf{d}}_{k,n}(\mathbf{X}, \Theta)) \leq \xi. \quad (37)$$

This implies that, with probability $1 - \rho$, for all n large enough and all $1 \leq j \leq k$, the j -th empirical cut $\hat{d}_{j,n}(\mathbf{X}, \Theta)$ is performed along the same coordinate as $d_j^*(\mathbf{X}, \Theta)$.

Now, for any cell A , since the regression function is not constant on A , one can find a theoretical cut d_A^* on A such that $L^*(d_A^*) > 0$. Thus, the cut d_A^* is made along an informative variable, in the sense that it is performed along one of the first S variables. Consequently, for all \mathbf{X}, Θ and for all $1 \leq j \leq k$, each theoretical cut $d_j^*(\mathbf{X}, \Theta)$ is made along one of the first S coordinates. The proof is then a consequence of inequality (37). \square

Acknowledgements

This work was supported by the European Research Council [SMAC-ERC-280032]. We greatly thank two referees for valuable comments and insightful suggestions.

References

- D. Amaratunga, J. Cabrera, and Y.-S. Lee. Enriched random forests. *Bioinformatics*, 24:2010–2014, 2008.
- Z.-H. Bai, L. Devroye, H.-K. Hwang, and T.-H. Tsai. Maxima in hypercubes. *Random Structures & Algorithms*, 27:290–309, 2005.
- O. Barndorff-Nielsen and M. Sobel. On the distribution of the number of admissible points in a vector random sample. *Theory of Probability & Its Applications*, 11:249–269, 1966.
- G. Biau. Analysis of a random forests model. *Journal of Machine Learning Research*, 13:1063–1095, 2012.
- G. Biau and L. Devroye. On the layered nearest neighbour estimate, the bagged nearest neighbour estimate and the random forest method in regression and classification. *Journal of Multivariate Analysis*, 101:2499–2518, 2010.

- G. Biau, L. Devroye, and G. Lugosi. Consistency of random forests and other averaging classifiers. *Journal of Machine Learning Research*, 9:2015–2033, 2008.
- S. Boucheron, G. Lugosi, and P. Massart. *Concentration inequalities: A nonasymptotic theory of independence*. Oxford University Press, 2013.
- L. Breiman. Bagging predictors. *Machine Learning*, 24:123–140, 1996.
- L. Breiman. Random forests. *Machine Learning*, 45:5–32, 2001.
- L. Breiman. *Consistency for a simple model of random forests*. Technical Report 670, UC Berkeley, 2004.
- L. Breiman, J. Friedman, R.A. Olshen, and C.J. Stone. *Classification and Regression Trees*. Chapman & Hall, New York, 1984.
- P. Bühlmann and B. Yu. Analyzing bagging. *The Annals of Statistics*, 30: 927–961, 2002.
- S. Cléménçon, M. Depecker, and N. Vayatis. Ranking forests. *Journal of Machine Learning Research*, 14:39–73, 2013.
- D.R. Cutler, T.C. Edwards Jr, K.H. Beard, A. Cutler, K.T. Hess, J. Gibson, and J.J. Lawler. Random forests for classification in ecology. *Ecology*, 88: 2783–2792, 2007.
- M. Denil, D. Matheson, and N. de Freitas. *Consistency of online random forests*, 2013. arXiv:1302.4853.
- L. Devroye, L. Györfi, and G. Lugosi. *A Probabilistic Theory of Pattern Recognition*. Springer, New York, 1996.
- R. Díaz-Uriarte and S. Alvarez de Andrés. Gene selection and classification of microarray data using random forest. *BMC Bioinformatics*, 7:1–13, 2006.
- R. Genuer. Variance reduction in purely random forests. *Journal of Non-parametric Statistics*, 24:543–562, 2012.
- P. Geurts, D. Ernst, and L. Wehenkel. Extremely randomized trees. *Machine Learning*, 63:3–42, 2006.
- L. Györfi, M. Kohler, A. Krzyżak, and H. Walk. *A Distribution-Free Theory of Nonparametric Regression*. Springer, New York, 2002.

- T. Hastie and R. Tibshirani. Generalized additive models. *Statistical Science*, 1:297–310, 1986.
- T. Hastie, R. Tibshirani, and J. Friedman. *The Elements of Statistical Learning. Second Edition*. Springer, New York, 2009.
- H. Ishwaran and U.B. Kogalur. Consistency of random survival forests. *Statistics & Probability Letters*, 80:1056–1064, 2010.
- H. Ishwaran, U.B. Kogalur, E.H. Blackstone, and M.S. Lauer. Random survival forests. *The Annals of Applied Statistics*, 2:841–860, 2008.
- A. Kleiner, A. Talwalkar, P. Sarkar, and M.I. Jordan. A scalable bootstrap for massive data. arXiv:1112.5016, 2012.
- B. Laurent and P. Massart. Adaptive estimation of a quadratic functional by model selection. *The Annals of Statistics*, 28:1302–1338, 2000.
- Y. Lin and Y. Jeon. Random forests and adaptive nearest neighbors. *Journal of the American Statistical Association*, 101:578–590, 2006.
- N. Meinshausen. Quantile regression forests. *Journal of Machine Learning Research*, 7:983–999, 2006.
- L. Mentch and G. Hooker. Ensemble trees and clts: Statistical inference for supervised learning. arXiv:1404.6473, 2014.
- A. Nobel. Histogram regression estimation using data-dependent partitions. *The Annals of Statistics*, 24:1084–1105, 1996.
- A.M. Prasad, L.R. Iverson, and A. Liaw. Newer classification and regression tree techniques: Bagging and random forests for ecological prediction. *Ecosystems*, 9:181–199, 2006.
- E. Scornet. On the asymptotics of random forests. arXiv:1409.2090, 2014.
- J. Shotton, A. Fitzgibbon, M. Cook, T. Sharp, M. Finocchio, R. Moore, A. Kipman, and A. Blake. Real-time human pose recognition in parts from single depth images. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1297–1304, 2011.
- C.J. Stone. Consistent nonparametric regression. *The Annals of Statistics*, 5:595–645, 1977.
- C.J. Stone. Additive regression and other nonparametric models. *The Annals of Statistics*, pages 689–705, 1985.

- V. Svetnik, A. Liaw, C. Tong, J.C. Culberson, R.P. Sheridan, and B.P. Feuston. Random forest: A classification and regression tool for compound classification and QSAR modeling. *Journal of Chemical Information and Computer Sciences*, 43:1947–1958, 2003.
- S. Wager. Asymptotic theory for random forests. arXiv:1405.0352, 2014.
- S. Wager, T. Hastie, and B. Efron. Standard errors for bagged predictors and random forests. arXiv:1311.4555, 2013.
- R. Zhu, D. Zeng, and M.R. Kosorok. *Reinforcement learning trees*. Technical Report, University of North Carolina, 2012.