

**SUPPLEMENTARY MATERIAL FOR THE PAPER
“MINIMAX OPTIMAL RATES FOR MONDRIAN TREES
AND FORESTS”**

BY JAOUAD MOURTADA^{*,†}, STÉPHANE GAÏFFAS^{*,‡} AND ERWAN
SCORNET^{*,†}

École polytechnique[†] and Université Paris Diderot[‡]

1. Introduction. This supplementary material to the paper “Minimax optimal rates for Mondrian trees and forests” gathers several proofs and technical details and definitions that were omitted in the main paper. Namely, we start with a glossary of notations, then give extra definitions and notations for trees and nested trees partitions in Section 2. Then, we provide proofs that were omitted in the main paper by order of appearance, namely the proofs of Proposition 2, Theorem 1, Proposition 3, Proposition 4 and Lemma 1.

Sign	Description
\mathcal{D}_n	Data set
μ	Distribution of X on $[0, 1]^d$
C , resp. $ C $	A generic cell $C \subset [0, 1]^d$, resp. half-perimeter of C
λ	Lifetime parameter of Mondrian process
$\text{MP}(\lambda, C)$	Distribution of a Mondrian process defined on cell C with lifetime parameter λ .
Π_λ , resp. $\Pi_\lambda C$	Partition drawn from $\text{MP}(\lambda, [0, 1]^d)$, resp. from $\text{MP}(\lambda, C)$
$C_\lambda(x)$	Cell of a Mondrian Tree with parameter λ containing x .
$D_\lambda(x)$	Diameter of $C_\lambda(x)$
K_λ	Number of cells in a Mondrian Tree partition Π_λ
$\hat{f}_{\lambda,n}^{(m)}(x)$	Mondrian Tree estimate at query point x based on the Mondrian partition $\Pi_\lambda^{(m)}$
$\hat{f}_{\lambda,n,M}(x)$	Mondrian Forest estimate at query point x based on the Mondrian partitions $\Pi_{\lambda,M} = (\Pi_\lambda^{(1)}, \dots, \Pi_\lambda^{(M)})$

*Data Science Initiative of École polytechnique

Sign	Description
$\bar{f}_\lambda^{(m)}(x)$	Expected value of the regression function f inside the cell $C_\lambda^{(m)}(x)$
$\tilde{f}_\lambda(x)$	Expected value of $\bar{f}_\lambda^{(m)}(x)$ over $\Pi_\lambda^{(m)} \sim \text{MP}(\lambda, [0, 1]^d)$
$\mathcal{N}(T), \mathcal{N}^\circ(T), \mathcal{L}(T)$	Nodes, interior nodes and leaves of a tree
$\Sigma = (\sigma_{\mathbf{v}})_{\mathbf{v} \in \mathcal{N}^\circ(T)}$	Set of splits for all nodes in the tree
$\sigma_{\mathbf{v}} = (j_{\mathbf{v}}, s_{\mathbf{v}})$	A split at node \mathbf{v} characterized by its split dimension $j_{\mathbf{v}} \in \{1, \dots, d\}$ and its threshold $s_{\mathbf{v}} \in [0, 1]$
$\tau_{\mathbf{v}}$	Birth time of a node \mathbf{v}

2. Specific notations. Let us now introduce some specific notations to describe the decision tree structure and the Mondrian Process.

2.1. *Trees and nested trees partitions.* A decision tree (T, Σ) is composed of the following components:

- A finite rooted ordered binary tree T , with nodes $\mathcal{N}(T)$, interior nodes $\mathcal{N}^\circ(T)$ and leaves $\mathcal{L}(T)$ (so that $\mathcal{N}(T)$ is the disjoint union of $\mathcal{N}^\circ(T)$ and $\mathcal{L}(T)$). The nodes $\mathbf{v} \in \mathcal{N}(T)$ are finite words on the alphabet $\{0, 1\}$, that is elements of the set $\{0, 1\}^* = \bigcup_{n \geq 0} \{0, 1\}^n$: the root ϵ of T is the empty word, and for every interior $\mathbf{v} \in \{0, 1\}^*$, its left child is $\mathbf{v}0$ (obtained by adding a 0 at the end of \mathbf{v}) while its right child is $\mathbf{v}1$ (obtained by adding a 1 at the end of \mathbf{v}).
- A family of *splits* $\Sigma = (\sigma_{\mathbf{v}})_{\mathbf{v} \in \mathcal{N}^\circ(T)}$ at each interior node, where each split $\sigma_{\mathbf{v}} = (j_{\mathbf{v}}, s_{\mathbf{v}})$ is characterized by its split dimension $j_{\mathbf{v}} \in \{1, \dots, d\}$ and its threshold $s_{\mathbf{v}} \in [0, 1]$.

We associate to $\Pi = (T, \Sigma)$ a partition $(C_{\mathbf{v}})_{\mathbf{v} \in \mathcal{L}(T)}$ of the unit cube $[0, 1]^d$, called a *tree partition* (or *guillotine partition*). For each node $\mathbf{v} \in \mathcal{N}(T)$, we define a hyper-rectangular region $C_{\mathbf{v}}$ recursively:

- The cell associated to the root of T is $[0, 1]^d$;
- For each $\mathbf{v} \in \mathcal{N}^\circ(T)$, we define

$$C_{\mathbf{v}0} := \{x \in C_{\mathbf{v}} : x_{j_{\mathbf{v}}} \leq s_{j_{\mathbf{v}}}\} \quad \text{and} \quad C_{\mathbf{v}1} := C_{\mathbf{v}} \setminus C_{\mathbf{v}0}.$$

The leaf cells $(C_{\mathbf{v}})_{\mathbf{v} \in \mathcal{L}(T)}$ form a partition of $[0, 1]^d$ by construction. In what follows, we will identify a tree with splits (T, Σ) with its associated tree partition, and a node $\mathbf{v} \in \mathcal{N}(T)$ with the cell $C_{\mathbf{v}} \subset [0, 1]^d$. The Mondrian process, described in the next Section, defines a distribution over nested tree partitions, defined below.

DEFINITION 1 (Nested tree partitions). A tree partition $\Pi' = (T', \Sigma')$ is a *refinement* of the tree partition $\Pi = (T, \Sigma)$ if T is a subtree of T' and, for every $\mathbf{v} \in \mathcal{N}(T) \subseteq \mathcal{N}(T')$, $\sigma_{\mathbf{v}} = \sigma'_{\mathbf{v}}$. A *nested tree partition* is a family $(\Pi_t)_{t \geq 0}$ of tree partitions such that, for every $t, t' \in \mathbf{R}^+$ with $t \leq t'$, $\Pi_{t'}$ is a refinement of Π_t . Such a family can be described as follows: let \mathbf{T} be the (in general infinite, and possibly complete) rooted binary tree, such that $\mathcal{N}(\mathbf{T}) = \bigcup_{t \geq 0} \mathcal{N}(T_t) \subseteq \{0, 1\}^*$. For each $\mathbf{v} \in \mathcal{N}(T)$, let $\tau_{\mathbf{v}} = \inf\{t \geq 0 \mid \mathbf{v} \in \mathcal{N}(T_t)\} < \infty$ denote the *birth time* of the node \mathbf{v} . Additionally, let $\sigma_{\mathbf{v}}$ be the value of the split $\sigma_{\mathbf{v}, t}$ in Π_t for $t > \tau_{\mathbf{v}}$ (which does not depend on t by the refinement property). Then, Π is completely characterized by \mathbf{T} , $\Sigma = (\sigma_{\mathbf{v}})_{\mathbf{v} \in \mathcal{N}(\mathbf{T})}$ and $\mathfrak{T} = (\tau_{\mathbf{v}})_{\mathbf{v} \in \mathcal{N}(\mathbf{T})}$.

2.2. *Mondrian Process.* To define rigorously the Mondrian Process, we introduce the function Φ_C , which maps any family of couples $(e_{\mathbf{v}}^j, u_{\mathbf{v}}^j) \in \mathbf{R}^+ \times [0, 1]$ indexed by the coordinates $j \in \{1, \dots, d\}$ and the nodes $\mathbf{v} \in \{0, 1\}^*$ to a nested tree partition $\Pi = \Phi_C((e_{\mathbf{v}}^j, u_{\mathbf{v}}^j)_{\mathbf{v}, j})$ of C . The splits $\sigma_{\mathbf{v}} = (j_{\mathbf{v}}, s_{\mathbf{v}})$ and birth times $\tau_{\mathbf{v}}$ of the nodes $\mathbf{v} \in \{0, 1\}^*$ are defined recursively, starting from the root ϵ :

- For the root node ϵ , we let $\tau_{\epsilon} = 0$ and $C_{\epsilon} = C$.
- At each node $\mathbf{v} \in \{0, 1\}^*$, given the labels of all its ancestors $\mathbf{v}' \sqsubset \mathbf{v}$ (so that in particular $\tau_{\mathbf{v}}$ and $C_{\mathbf{v}}$ are determined), denote $C_{\mathbf{v}} = \prod_{j=1}^d [a_{\mathbf{v}}^j, b_{\mathbf{v}}^j]$. Then, select the split dimension $j_{\mathbf{v}} \in \{1, \dots, d\}$ and its location $s_{\mathbf{v}}$ as follows:

$$(2.1) \quad j_{\mathbf{v}} = \operatorname{argmin}_{j=1, \dots, d} \frac{e_{\mathbf{v}}^j}{b_{\mathbf{v}}^j - a_{\mathbf{v}}^j}, \quad s_{\mathbf{v}} = a_{\mathbf{v}}^{j_{\mathbf{v}}} + (b_{\mathbf{v}}^{j_{\mathbf{v}}} - a_{\mathbf{v}}^{j_{\mathbf{v}}}) \cdot u_{\mathbf{v}}^{j_{\mathbf{v}}},$$

where we break ties in the choice of $j_{\mathbf{v}}$ e.g., by choosing the smallest index j in the argmin. The node \mathbf{v} is then split at time $\tau_{\mathbf{v}} + e_{\mathbf{v}}^{j_{\mathbf{v}}}/(b_{\mathbf{v}}^{j_{\mathbf{v}}} - a_{\mathbf{v}}^{j_{\mathbf{v}}}) = \tau_{\mathbf{v}0} = \tau_{\mathbf{v}1}$, we let $C_{\mathbf{v}0} = \{x \in C_{\mathbf{v}} : x_{j_{\mathbf{v}}} \leq s_{\mathbf{v}}\}$, $C_{\mathbf{v}1} = C_{\mathbf{v}} \setminus C_{\mathbf{v}0}$ and recursively apply the procedure to its children $\mathbf{v}0$ and $\mathbf{v}1$.

For each $\lambda \in \mathbf{R}^+$, the tree partition $\Pi_{\lambda} = \Phi_{\lambda, C}((e_{\mathbf{v}}^j, u_{\mathbf{v}}^j)_{\mathbf{v}, j})$ is the *pruning of Π at time λ* , obtained by removing all the splits in Π that occurred strictly after λ , so that the leaves of the tree are the maximal nodes (in the prefix order) \mathbf{v} such that $\tau_{\mathbf{v}} \leq \lambda$.

DEFINITION 2 (Mondrian process). Let $(E_{\mathbf{v}}^j, U_{\mathbf{v}}^j)_{\mathbf{v}, j}$ be a family of independent random variables, with $E_{\mathbf{v}}^j \sim \operatorname{Exp}(1)$, $U_{\mathbf{v}}^j \sim \mathcal{U}([0, 1])$. The *Mondrian process* $\operatorname{MP}(C)$ on C is the distribution of the random nested tree partition $\Phi_C((E_{\mathbf{v}}^j, U_{\mathbf{v}}^j)_{\mathbf{v}, j})$. In addition, we denote $\operatorname{MP}(\lambda, C)$ the distribution of $\Phi_{\lambda, C}((E_{\mathbf{v}}^j, U_{\mathbf{v}}^j)_{\mathbf{v}, j})$.

3. Proof of Proposition 2. At a high level, the idea of the proof is to modify the construction of the Mondrian partition (and hence, the distribution of the underlying process) without affecting the expected number of cells. More precisely, we show a recursive way to transform the Mondrian process that leaves $\mathbb{E}[K_\lambda]$ unchanged, and which eventually leads to a random partition $\tilde{\Pi}_\lambda$ for which this quantity can be computed directly and equals $(1 + \lambda)^d$. We will in fact show the result for a general box C (not just the unit cube). The proof proceeds in two steps:

1. Define a modified process $\tilde{\Pi}$, and show that $\mathbb{E}[\tilde{K}_\lambda] = \prod_{j=1}^d (1 + \lambda |C^j|)$.
2. It remains to show that $\mathbb{E}[K_\lambda] = \mathbb{E}[\tilde{K}_\lambda]$. For this, it is sufficient to show that the distribution of the birth times $\tau_{\mathbf{v}}$ and $\tilde{\tau}_{\mathbf{v}}$ of the node \mathbf{v} is the same for both processes. This is done by induction on \mathbf{v} , by showing that the splits at one node of both processes have the same conditional distribution given the splits at previous nodes.

Let $(E_{\mathbf{v}}^j, U_{\mathbf{v}}^j)_{\mathbf{v} \in \{0,1\}^*, 1 \leq j \leq d}$ be a family of independent random variables with $E_{\mathbf{v}}^j \sim \text{Exp}(1)$ and $U_{\mathbf{v}}^j \sim \mathcal{U}([0, 1])$. By definition, $\Pi = \Phi_C((E_{\mathbf{v}}^j, U_{\mathbf{v}}^j)_{\mathbf{v}, j})$ (Φ_C being defined in Section 3) follows a Mondrian process distribution $\text{MP}(C)$. Denote for every node $\mathbf{v} \in \{0, 1\}^*$ $C_{\mathbf{v}}$ the cell of \mathbf{v} , $\tau_{\mathbf{v}}$ its birth time, as well as its split time $T_{\mathbf{v}}$, dimension $J_{\mathbf{v}}$, and threshold $S_{\mathbf{v}}$ (note that $T_{\mathbf{v}} = \tau_{\mathbf{v}0} = \tau_{\mathbf{v}1}$). In addition, for $\lambda \in \mathbf{R}^+$, denote $\Pi_\lambda \sim \text{MP}(\lambda, C)$ the tree partition restricted to time λ , and $K_\lambda \in \mathbf{N} \cup \{+\infty\}$ its number of nodes.

Construction of the modified process. Now, consider the following modified nested partition of C , denoted $\tilde{\Pi}$, and defined through its split times, dimension and threshold $\tilde{T}_{\mathbf{v}}, \tilde{J}_{\mathbf{v}}, \tilde{S}_{\mathbf{v}}$ (which determine the birth times $\tau_{\mathbf{v}}$ and cells $C_{\mathbf{v}}$), and *current j -dimensional node* $\mathbf{v}_j(\mathbf{v}) \in \{0, 1\}^*$ ($1 \leq j \leq d$) at each node \mathbf{v} . First, for every $j = 1, \dots, d$, let $\Pi^{j'} = \Phi_{C^j}((E_{\mathbf{v}}^j, U_{\mathbf{v}}^j)_{\mathbf{v} \in \{0,1\}^*}) \sim \text{MP}(C^j)$ be the nested partition of the interval C^j determined by $(E_{\mathbf{v}}^j, U_{\mathbf{v}}^j)_{\mathbf{v}}$; its split times and thresholds are denoted $(S_{\mathbf{v}}^{j'}, T_{\mathbf{v}}^{j'})$. Then, $\tilde{\Pi}$ is defined recursively as follows:

- At the root node ϵ , let $\tilde{\tau}_\epsilon = 0$, $\tilde{C}_\epsilon = C$ and $\mathbf{v}_j(\epsilon) := \epsilon$ for $1 \leq j \leq d$.
- At node \mathbf{v} , given $(\tau_{\mathbf{v}'}, C_{\mathbf{v}'}, \mathbf{v}_j(\mathbf{v}'))_{\mathbf{v}' \sqsubseteq \mathbf{v}}$ (i.e., given $(\tilde{J}_{\mathbf{v}'}, \tilde{S}_{\mathbf{v}'}, \tilde{T}_{\mathbf{v}'})_{\mathbf{v}' \sqsubseteq \mathbf{v}}$) define:

$$(3.1) \quad \tilde{T}_{\mathbf{v}} = \min_{1 \leq j \leq d} T_{\mathbf{v}_j(\mathbf{v})}^{j'}, \quad \tilde{J}_{\mathbf{v}} := \operatorname{argmin}_{1 \leq j \leq d} T_{\mathbf{v}_j(\mathbf{v})}^{j'}, \quad \tilde{S}_{\mathbf{v}} = S_{\mathbf{v}_j(\mathbf{v})}^{j'},$$

$$(3.2) \quad \mathbf{v}_j(\mathbf{v}a) = \begin{cases} \mathbf{v}_j(\mathbf{v})a & \text{if } j = \tilde{J}_{\mathbf{v}} \\ \mathbf{v}_j(\mathbf{v}) & \text{else.} \end{cases}$$

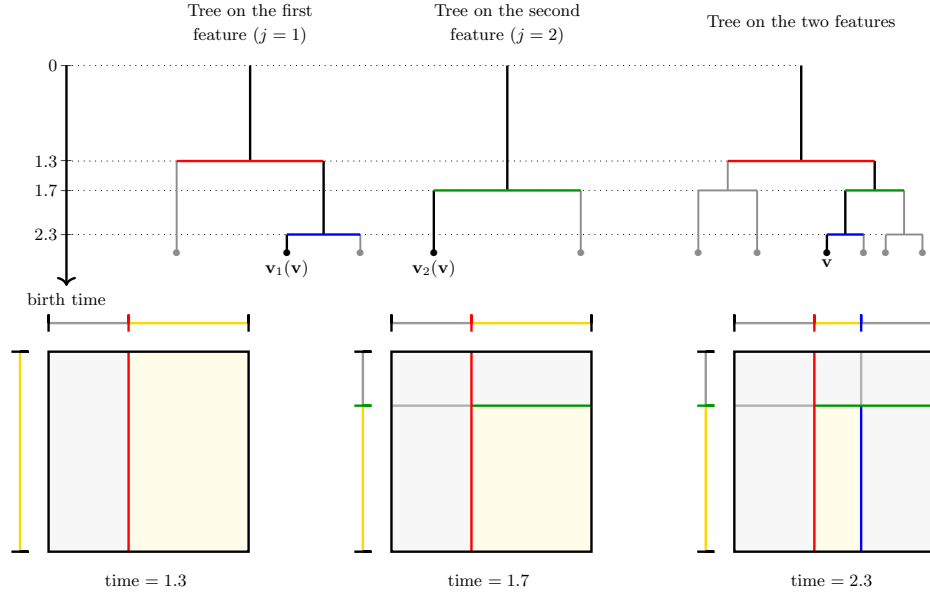


FIG 1. *Modified construction in dimension two. At the top, from left to right: trees associated to partitions Π^1, Π^2 and $\tilde{\Pi}$ respectively. At the bottom, from left to right: successive splits in $\tilde{\Pi}$ leading to the leaf \mathbf{v} (depicted in yellow).*

Finally, for every $\lambda \in \mathbf{R}^+$, define $\tilde{\Pi}_\lambda$ and \tilde{K}_λ as before from $\tilde{\Pi}$. This construction is illustrated in Figure 1.

Computation of $\mathbb{E}[\tilde{K}_\lambda]$. Now, it can be seen that the partition $\tilde{\Pi}_\lambda$ is a rectangular grid which is the “product” of the partitions Π^{t_j} of the intervals C^j , $1 \leq j \leq d$. Indeed, let $x \in [0, 1]^d$, and let $\tilde{C}_\lambda(x)$ be the cell in $\tilde{\Pi}_\lambda$ that contains x ; we need to show that $\tilde{C}_\lambda(x) = \prod_{j=1}^d C_\lambda^{t_j}(x)$, where $C_\lambda^{t_j}(x)$ is the subinterval of C^j in the partition Π^{t_j} that contains x_j . The proof proceeds in several steps:

- First, Equation (3.1) shows that, for every node \mathbf{v} , we have $\tilde{C}_\mathbf{v} = \prod_{1 \leq j \leq d} C_{\mathbf{v}_j(\mathbf{v})}^{t_j}$, since the successive splits on the j -th coordinate of $\tilde{C}_\mathbf{v}$ are precisely the ones of $C_{\mathbf{v}_j(\mathbf{v})}^{t_j}$.
- Second, it follows from (3.1) that $\tilde{T}_\mathbf{v} = \min_{1 \leq j \leq d} T_{\mathbf{v}_j(\mathbf{v})}^{t_j}$; also, since the cell $C_\mathbf{v}$ is formed when its last split is performed, $\tilde{\tau}_\mathbf{v} = \max_{1 \leq j \leq d} \tau_{\mathbf{v}_j(\mathbf{v})}^{t_j}$.
- Let $\tilde{\mathbf{v}}$ be the node such that $\tilde{C}_{\tilde{\mathbf{v}}} = \tilde{C}_\lambda(x)$, and \mathbf{v}^j be such that $C_{\mathbf{v}^j}^{t_j} = C_\lambda^{t_j}(x_j)$. By the first point, it suffices to show that $\mathbf{v}_j(\tilde{\mathbf{v}}) = \mathbf{v}_j^j$ for $1 \leq j \leq d$.

- Observe that $\tilde{\mathbf{v}}$ (resp. \mathbf{v}'_j) is characterized by the fact that $x \in \tilde{C}_{\tilde{\mathbf{v}}}$ and $\tilde{\tau}_{\tilde{\mathbf{v}}} \leq \lambda < \tilde{T}_{\tilde{\mathbf{v}}}$ (resp. $x_j \in C'_{\mathbf{v}'_j}$ and $\tau'_{\mathbf{v}'_j} \leq \lambda < T'_{\mathbf{v}'_j}$). But since $\tilde{C}_{\tilde{\mathbf{v}}} = \prod_{1 \leq j \leq d} C'_{\mathbf{v}'_j}$ (first point), $x \in \tilde{C}_{\tilde{\mathbf{v}}}$ implies $x_j \in C'_{\mathbf{v}'_j}$. Likewise, since $\tilde{\tau}_{\tilde{\mathbf{v}}} = \max_{1 \leq j \leq d} \tau'_{\mathbf{v}'_j}$ and $\tilde{T}_{\tilde{\mathbf{v}}} = \min_{1 \leq j \leq d} T'_{\mathbf{v}'_j}$ (second point), $\tilde{\tau}_{\tilde{\mathbf{v}}} \leq \lambda < \tilde{T}_{\tilde{\mathbf{v}}}$ implies $\tau'_{\mathbf{v}'_j} \leq \lambda < T'_{\mathbf{v}'_j}$. Since these properties characterize \mathbf{v}'_j , we have $\mathbf{v}_j(\tilde{\mathbf{v}}) = \mathbf{v}'_j$, which concludes the proof.

Hence, the partition $\tilde{\Pi}_\lambda$ is the product of the partitions $\Pi^j = \Phi_{C^j}((E_{\mathbf{v}}^j, U_{\mathbf{v}}^j)_{\mathbf{v}})_\lambda$ of the intervals C^j , $1 \leq j \leq d$, which are independent Mondrians distributed as $\text{MP}(\lambda, C^j)$. By Fact 1, the splits of the Mondrian partition $\text{MP}(\lambda, C^j)$ are distributed as a Poisson point process on C^j of intensity λ , so that the expected number of cells in such a partition is $1 + \lambda|C^j|$. Since $\tilde{\Pi}_\lambda$ is a “product” of such independent partitions, we have:

$$(3.3) \quad \mathbb{E}[\tilde{K}_\lambda] = \prod_{j=1}^d (1 + \lambda|C^j|).$$

Equality of $\mathbb{E}[K_\lambda]$ and $\mathbb{E}[\tilde{K}_\lambda]$. In order to establish Proposition 2, it is thus sufficient to prove that $\mathbb{E}[K_\lambda] = \mathbb{E}[\tilde{K}_\lambda]$. First, note that, since the number of cells in a partition is one plus the number of splits (each split increases the number of cells by one)

$$K_\lambda = 1 + \sum_{\mathbf{v} \in \{0,1\}^*} \mathbf{1}(T_{\mathbf{v}} \leq \lambda)$$

so that we have, respectively,

$$(3.4) \quad \mathbb{E}[K_\lambda] = 1 + \sum_{\mathbf{v} \in \{0,1\}^*} \mathbb{P}(T_{\mathbf{v}} \leq \lambda)$$

$$(3.5) \quad \mathbb{E}[\tilde{K}_\lambda] = 1 + \sum_{\mathbf{v} \in \{0,1\}^*} \mathbb{P}(\tilde{T}_{\mathbf{v}} \leq \lambda).$$

Hence, it suffices to show that $\mathbb{P}(T_{\mathbf{v}} \leq \lambda) = \mathbb{P}(\tilde{T}_{\mathbf{v}} \leq \lambda)$ for every $\mathbf{v} \in \{0,1\}^*$ and $\lambda \geq 0$, *i.e.* that $T_{\mathbf{v}}$ and $\tilde{T}_{\mathbf{v}}$ have the same distribution for every \mathbf{v} .

In order to establish this, we show that, for every $\mathbf{v} \in \{0,1\}^*$, the conditional distribution of $(\tilde{T}_{\mathbf{v}}, \tilde{J}_{\mathbf{v}}, \tilde{S}_{\mathbf{v}})$ given $\tilde{\mathcal{F}}_{\mathbf{v}} = \sigma((\tilde{T}_{\mathbf{v}'}, \tilde{J}_{\mathbf{v}'}, \tilde{S}_{\mathbf{v}'})_{\mathbf{v}' \sqsubset \mathbf{v}})$ has the same form as the conditional distribution of $(T_{\mathbf{v}}, J_{\mathbf{v}}, S_{\mathbf{v}})$ given $\mathcal{F}_{\mathbf{v}} = \sigma((T_{\mathbf{v}'}, J_{\mathbf{v}'}, S_{\mathbf{v}'})_{\mathbf{v}' \sqsubset \mathbf{v}})$, in the sense that there exists a family of conditional distributions $(\Psi_{\mathbf{v}})_{\mathbf{v}}$ such that, for every \mathbf{v} , the conditional distribution of

$(T_{\mathbf{v}}, J_{\mathbf{v}}, S_{\mathbf{v}})$ given $\mathcal{F}_{\mathbf{v}}$ is $\Psi_{\mathbf{v}}(\cdot | (T_{\mathbf{v}'}, J_{\mathbf{v}'}, S_{\mathbf{v}'}) , \mathbf{v}' \sqsubset \mathbf{v})$ and the conditional distribution of $(\tilde{T}_{\mathbf{v}}, \tilde{J}_{\mathbf{v}}, \tilde{S}_{\mathbf{v}})$ given $\tilde{\mathcal{F}}_{\mathbf{v}}$ is $\Psi_{\mathbf{v}}(\cdot | (\tilde{T}_{\mathbf{v}'}, \tilde{J}_{\mathbf{v}'}, \tilde{S}_{\mathbf{v}'}) , \mathbf{v}' \sqsubset \mathbf{v})$.

First, recall that the variables $(E_{\mathbf{v}'}^j, U_{\mathbf{v}'}^j)_{\mathbf{v}' \in \{0,1\}^*, 1 \leq j \leq d}$ are independent, so $(E_{\mathbf{v}}^j, U_{\mathbf{v}}^j)_{1 \leq j \leq d}$ is independent from $\mathcal{F}_{\mathbf{v}}$. Hence, conditionally on $\mathcal{F}_{\mathbf{v}}$, $E_{\mathbf{v}}^j, U_{\mathbf{v}}^j$, $1 \leq j \leq d$ are independent with $E_{\mathbf{v}}^j \sim \text{Exp}(1)$ and $U_{\mathbf{v}}^j \sim \mathcal{U}([0, 1])$. Also, recall that if T_1, \dots, T_d are independent exponential random variables of intensities $\lambda_1, \dots, \lambda_d$, and if $T = \min_{1 \leq j \leq d} T_j$ and $J = \text{argmin}_{1 \leq j \leq d} T_j$, then $\mathbb{P}(J = j) = \lambda_j / \sum_{j'=1}^d \lambda_{j'}$, $T \sim \text{Exp}(\sum_{j=1}^d \lambda_j)$ and J and T are independent. Hence, conditionally on $\mathcal{F}_{\mathbf{v}}$, $T_{\mathbf{v}} - \tau_{\mathbf{v}} = \min_{1 \leq j \leq d} E_{\mathbf{v}}^j / |C_{\mathbf{v}}^j| \sim \text{Exp}(\sum_{j=1}^d |C_{\mathbf{v}}^j|) = \text{Exp}(|C_{\mathbf{v}}|)$, $J_{\mathbf{v}} := \text{argmin}_{1 \leq j \leq d} E_{\mathbf{v}}^j / |C_{\mathbf{v}}^j|$ equals j with probability $|C_{\mathbf{v}}^j| / |C_{\mathbf{v}}|$, $T_{\mathbf{v}}, J_{\mathbf{v}}$ are independent and $(S_{\mathbf{v}} | T_{\mathbf{v}}, J_{\mathbf{v}}) \sim \mathcal{U}(C_{\mathbf{v}}^{J_{\mathbf{v}}})$.

Now consider the conditional distribution of $(\tilde{T}_{\mathbf{v}}, \tilde{J}_{\mathbf{v}}, \tilde{S}_{\mathbf{v}})$ given $\tilde{\mathcal{F}}_{\mathbf{v}}$. Let $(\mathbf{v}_v)_{v \in \mathbf{N}}$ be a path in $\{0, 1\}^*$ from the root: $\mathbf{v}_0 := \epsilon$, \mathbf{v}_{v+1} is a child of \mathbf{v}_v for $v \in \mathbf{N}$, and $\mathbf{v}_v \sqsubseteq \mathbf{v}$ for $0 \leq v \leq \text{depth}(\mathbf{v})$. Define for $v \in \mathbf{N}$, $E_v^j = E_{\mathbf{v}_v}^j$ and $U_v^j = U_{\mathbf{v}_v}^j$ if \mathbf{v}_{v+1} is the left child of \mathbf{v}_v , and $1 - U_{\mathbf{v}_v}^j$ otherwise. Then, the variables $(E_v^j, U_v^j)_{v \in \mathbf{N}, 1 \leq j \leq d}$ are independent, with $E_v^j \sim \text{Exp}(1)$, $U_v^j \sim \mathcal{U}([0, 1])$, so that the following Lemma applies.

LEMMA 1. *Let $(E_v^j, U_v^j)_{v \in \mathbf{N}^*, 1 \leq j \leq d}$ be a family of independent random variables, with $U_v^j \sim \mathcal{U}([0, 1])$ and $E_v^j \sim \text{Exp}(1)$. Let $a_1, \dots, a_d > 0$. For $1 \leq j \leq d$, define the sequence $(T_v^j, L_v^j)_{v \in \mathbf{N}}$ as follows:*

- $L_0^j = a_j$, $T_0^j = \frac{E_0^j}{a_j}$;
- for $v \in \mathbf{N}$, $L_{v+1}^j = U_v^j L_v^j$, $T_{v+1}^j = T_v^j + \frac{E_{v+1}^j}{L_{v+1}^j}$.

Define recursively the variables \tilde{V}_v^j ($v \in \mathbf{N}, 1 \leq j \leq d$) as well as $\tilde{J}_v, \tilde{T}_v, \tilde{U}_v$ ($v \in \mathbf{N}$) as follows:

- $\tilde{V}_0^j = 0$ for $j = 1, \dots, d$.
- for $v \in \mathbf{N}$, given \tilde{V}_v^j ($1 \leq j \leq d$), denoting $\tilde{T}_v^j = T_{\tilde{V}_v^j}^j$ and $\tilde{U}_v^j = U_{\tilde{V}_v^j}^j$, set

$$(3.6) \quad \begin{aligned} \tilde{J}_v &= \text{argmin}_{1 \leq j \leq d} \tilde{T}_v^j, & \tilde{T}_v &= \min_{1 \leq j \leq d} \tilde{T}_v^j = \tilde{T}_v^{\tilde{J}_v}, & \tilde{U}_v &= \tilde{U}_v^{\tilde{J}_v}, \\ & & \text{and } \tilde{V}_{v+1}^j &= \tilde{V}_v^j + \mathbf{1}(\tilde{J}_v = j). \end{aligned}$$

Then, the conditional distribution of $(\tilde{J}_v, \tilde{T}_v, \tilde{U}_v)$ given $\mathcal{F}_v = \sigma((\tilde{J}_{v'}, \tilde{T}_{v'}, \tilde{U}_{v'}), 0 \leq v' < v)$ is the following (denoting $\tilde{L}_v^j = L_{\tilde{V}_v^j}^j$):

- $\tilde{J}_v, \tilde{T}_v, \tilde{U}_v$ are independent,

- $\mathbb{P}(\tilde{J}_v = j \mid \mathcal{F}_v) = \tilde{L}_v^j / (\sum_{j'=1}^d \tilde{L}_v^{j'})$,
- $\tilde{T}_v - \tilde{T}_{v-1} \sim \text{Exp}(\sum_{j=1}^d \tilde{L}_v^j)$ (with the convention $\tilde{T}_{-1} = 0$) and $\tilde{U}_v \sim \mathcal{U}([0, 1])$.

In addition, note that, with the notations of Lemma 1, a simple induction shows that $\tilde{J}_v = \tilde{J}_{\mathbf{v}_v}$, $\tilde{T}_v = \tilde{T}_{\mathbf{v}_v}$, $\tilde{U}_v = \tilde{U}_{\mathbf{v}_v}$ and $L_v^j = |\tilde{C}_{\mathbf{v}_v}^j|$, so that $\mathcal{F}_v = \mathcal{F}_{\mathbf{v}_v}$. Applying Lemma 1 for $v = \text{depth}(\mathbf{v})$ (so that $\mathbf{v}_v = \mathbf{v}$) therefore gives the following: conditionally on $\mathcal{F}_{\mathbf{v}}$, the variables $\tilde{T}_{\mathbf{v}}, \tilde{J}_{\mathbf{v}}, \tilde{U}_{\mathbf{v}}$ are independent, $\tilde{T}_{\mathbf{v}} - \tilde{\tau}_{\mathbf{v}} \sim \text{Exp}(|\tilde{C}_{\mathbf{v}}^j|)$, $\mathbb{P}(\tilde{J}_{\mathbf{v}} = j \mid \mathcal{F}_{\mathbf{v}}) = |\tilde{C}_{\mathbf{v}}^j| / (\sum_{j'=1}^d |\tilde{C}_{\mathbf{v}}^{j'}|)$ and $\tilde{U}_{\mathbf{v}} \sim \mathcal{U}([0, 1])$, so that $(\tilde{S}_{\mathbf{v}} \mid \mathcal{F}_{\mathbf{v}}, \tilde{T}_{\mathbf{v}}, \tilde{J}_{\mathbf{v}}) \sim \mathcal{U}(\tilde{C}_{\mathbf{v}}^{\tilde{J}_{\mathbf{v}}})$. Hence, we have proven that, for every \mathbf{v} , the conditional distribution of $(T_{\mathbf{v}}, J_{\mathbf{v}}, S_{\mathbf{v}})$ given $\mathcal{F}_{\mathbf{v}}$ is the same as that of $(\tilde{T}_{\mathbf{v}}, \tilde{J}_{\mathbf{v}}, \tilde{S}_{\mathbf{v}})$ given $\tilde{\mathcal{F}}_{\mathbf{v}}$. By induction on \mathbf{v} , since $\mathcal{F}_{\epsilon} = \tilde{\mathcal{F}}_{\epsilon}$ is the trivial σ -algebra, this shows that $T_{\mathbf{v}}$ and $\tilde{T}_{\mathbf{v}}$ have the same distribution for every \mathbf{v} . Plugging this into (3.4) and (3.5) and combining it with (3.3) completes the proof of Proposition 2. \square

PROOF OF LEMMA 1. We show by induction on $v \in \mathbb{N}$ the following property: conditionally on $\mathcal{F}_{\mathbf{v}}$, $(\tilde{T}_v^j, \tilde{U}_v^j)_{1 \leq j \leq d}$ are independent, $\tilde{T}_v^j - \tilde{T}_{v-1} \sim \text{Exp}(L_v^j)$ and $\tilde{U}_v^j \sim \mathcal{U}([0, 1])$.

Initialization For $v = 0$ (with \mathcal{F}_0 the trivial σ -algebra), since $\tilde{V}_0^j = 0$ we have $\tilde{T}_0^j = E_0^j / a_j \sim \text{Exp}(a_j) = \text{Exp}(L_0^j)$, $\tilde{U}_0^j = U_0^j \sim \mathcal{U}([0, 1])$ and these random variables are independent.

Inductive step Let $v \in \mathbb{N}$, and assume the property is true up to step v . Conditionally on \mathcal{F}_{v+1} , *i.e.* on $\mathcal{F}_v, \tilde{T}_v, \tilde{J}_v, \tilde{U}_v$, we have:

- for $j \neq \tilde{J}_v$, the variables $\tilde{T}_{v+1}^j - \tilde{T}_{v-1} = \tilde{T}_v^j - \tilde{T}_{v-1}$ are independent $\text{Exp}(\tilde{L}_v^j) = \text{Exp}(\tilde{L}_{v+1}^j)$ random variables (when conditioned only on \mathcal{F}_v , by the induction hypothesis), conditioned on $\tilde{T}_{v+1}^j - \tilde{T}_{v-1} \geq \tilde{T}_v - \tilde{T}_{v-1}$, so by the memory-less property of exponential random variables $\tilde{T}_{v+1}^j - \tilde{T}_v = (\tilde{T}_{v+1}^j - \tilde{T}_{v-1}) - (\tilde{T}_v - \tilde{T}_{v-1}) \sim \text{Exp}(\tilde{L}_{v+1}^j)$ (and those variables are independent).
- for $j \neq \tilde{J}_v$, the variables $\tilde{U}_{v+1}^j = \tilde{U}_v^j$ are independent $\mathcal{U}([0, 1])$ random variables (conditionally on \mathcal{F}_v), conditioned on the independent variables $\tilde{T}_v, \tilde{J}_v, \tilde{U}_v$, so they remain independent $\mathcal{U}([0, 1])$ random variables.
- $(\tilde{T}_{v+1}^{\tilde{J}_v} - \tilde{T}_v, \tilde{U}_{v+1}^{\tilde{J}_v}) = (E_{\tilde{V}_{v+1}^{\tilde{J}_v}}^{\tilde{J}_v} / \tilde{L}_{v+1}^{\tilde{J}_v}, U_{\tilde{V}_{v+1}^{\tilde{J}_v}}^{\tilde{J}_v})$ is distributed, conditionally on \mathcal{F}_{v+1} , *i.e.* on $\tilde{J}_v, \tilde{T}_v, \tilde{V}_{v+1}^{\tilde{J}_v}, \tilde{L}_{v+1}^{\tilde{J}_v}$, as $\text{Exp}(\tilde{L}_{v+1}^{\tilde{J}_v}) \otimes \mathcal{U}([0, 1])$, and independent of $(\tilde{T}_{v+1}^j, \tilde{U}_{v+1}^j)_{j \neq \tilde{J}_v}$.

This completes the proof by induction.

Let $v \in \mathbf{N}$. We have established that, conditionally on \mathcal{F}_v , the variables $(\tilde{T}_v^j, \tilde{U}_v^j)_{1 \leq j \leq d}$ are independent, with $\tilde{T}_v^j - \tilde{T}_{v-1} \sim \text{Exp}(\tilde{L}_v^j)$ and $\tilde{U}_v^j \sim \mathcal{U}([0, 1])$. In particular, conditionally on \mathcal{F}_v , \tilde{U}_v is independent from $(\tilde{J}_v, \tilde{T}_v)$, $\tilde{U}_v \sim \mathcal{U}([0, 1])$, and (by the property of the minimum of independent exponential random variables) \tilde{J}_v is independent of \tilde{T}_v , $\tilde{T}_v \sim \text{Exp}(\sum_{j=1}^d \tilde{L}_v^j)$ and $\mathbb{P}(\tilde{J}_v = j | \mathcal{F}_v) = \tilde{L}_v^j / (\sum_{j'=1}^d \tilde{L}_v^{j'})$. This concludes the proof of Lemma 1. \square

4. Proof of Theorem 1. Recall that a Mondrian Forest estimate with lifetime parameter λ , is defined, for all $x \in [0, 1]^d$, by

$$\hat{f}_{\lambda, n, M}(x) = \hat{f}_{\lambda, n, M}(x, \Pi_{\lambda, M}) = \frac{1}{M} \sum_{m=1}^M \hat{f}_{\lambda, n}^{(m)}(x, \Pi_{\lambda}^{(m)}),$$

where $\hat{f}_{\lambda, n}^{(m)}(x, \Pi_{\lambda}^{(m)})$ denotes the Mondrian Tree based on the random partition $\Pi_{\lambda}^{(m)}$ and $\Pi_{\lambda, M} = (\Pi_{\lambda}^{(1)}, \dots, \Pi_{\lambda}^{(M)})$. To ease notation, we will write $\hat{f}_{\lambda, n}^{(m)}(x)$ instead of $\hat{f}_{\lambda, n}^{(m)}(x, \Pi_{\lambda}^{(m)})$. First, note that, by Jensen's inequality,

$$\begin{aligned} R(\hat{f}_{\lambda, n, M}) &= \mathbb{E}_{(X, \Pi_{\lambda, M})} [(\hat{f}_{\lambda, n, M}(x, \Pi_{\lambda, M}) - f(X))^2] \\ &\leq \frac{1}{M} \sum_{m=1}^M \mathbb{E}_{(X, \Pi_{\lambda}^{(m)})} [(\hat{f}_{\lambda, n}^{(m)}(X, \Pi_{\lambda}^{(m)}) - f(X))^2] \\ &\leq \mathbb{E}_{(X, \Pi_{\lambda}^{(1)})} [(\hat{f}_{\lambda, n}^{(1)}(X, \Pi_{\lambda}^{(1)}) - f(X))^2], \end{aligned}$$

since each Mondrian tree has the same distribution. Therefore, it is sufficient to prove that a single Mondrian tree is consistent. Now, since Mondrian partitions are independent of the dataset \mathcal{D}_n , we can apply Theorem 4.2 from [3], which states that a Mondrian tree estimate is consistent if

- (i) $D_{\lambda}(X) \rightarrow 0$ in probability, as $n \rightarrow \infty$,
- (ii) $K_{\lambda}/n \rightarrow \infty$ in probability, as $n \rightarrow \infty$,

where $D_{\lambda}(X)$ is the diameter of the cell of the Mondrian tree that contains X , and K_{λ} is the number of cells in the Mondrian tree. Note that the initial assumptions in Theorem 4.2 in [3] contains deterministic convergence, but can be relaxed to convergences in probability by a close inspection of the proof. In the sequel, we prove that an individual Mondrian tree satisfies (i) and (ii) which will conclude the proof. To prove (i), just note that, according to Corollary 1,

$$\mathbb{E}[D_{\lambda}(X)^2] = \mathbb{E}[\mathbb{E}[D_{\lambda}(X)^2 | X]] \leq \frac{4d}{\lambda^2},$$

which tends to zero, since $\lambda = \lambda_n \rightarrow \infty$, as $n \rightarrow \infty$. Thus, condition (i) holds. Now, to prove (ii), observe that

$$\mathbb{E}\left[\frac{K_\lambda}{n}\right] = \frac{(1 + \lambda)^d}{n},$$

which tends to zero since $\lambda_n^d/n \rightarrow 0$ by assumption, as $n \rightarrow \infty$. This concludes the proof of Theorem 1. \square

5. Proof of Proposition 3. Let $\Pi_\lambda^{(1)}$ be the Mondrian partition of $[0, 1]$ used to construct the randomized estimator $\tilde{f}_{\lambda,n}^{(1)}$. Denote by $\bar{f}_\lambda^{(1)}$ the random function $\bar{f}_\lambda^{(1)}(x) = \mathbb{E}_X[f(X) | X \in C_\lambda(x)]$, and define $\tilde{f}_\lambda(x) = \mathbb{E}\left[\bar{f}_\lambda^{(1)}(x)\right]$ (which is deterministic). For the seek of clarity, we will drop the exponent “(1)” in all notations, keeping in mind that we consider only one particular Mondrian partition, whose associated Mondrian Tree estimate is denoted by $\hat{f}_{\lambda,n}$. Recall the bias-variance decomposition (7.4) for Mondrian trees:

$$(5.1) \quad R(\hat{f}_{\lambda,n}^{(1)}) = \mathbb{E}[(f(X) - \bar{f}_\lambda(X))^2] + \mathbb{E}[(\bar{f}_\lambda(X) - \tilde{f}_{\lambda,n}^{(1)}(X))^2].$$

We will provide lower bounds for the first term (the bias, depending on λ) and the second (the variance, depending on both λ and n), which will lead to the stated lower bound on the risk, valid for every value of λ .

Lower bound on the bias. As we will see, the point-wise bias $\mathbb{E}[(\bar{f}_\lambda(x) - f(x))^2]$ can be computed explicitly given our assumptions. Let $x \in [0, 1]$. Since $\tilde{f}_\lambda(x) = \mathbb{E}[\bar{f}_\lambda(x)]$, we have

$$(5.2) \quad \mathbb{E}[(\bar{f}_\lambda(x) - f(x))^2] = \text{Var}(\bar{f}_\lambda(x)) + (\tilde{f}_\lambda(x) - f(x))^2.$$

By Proposition 1, the cell of x in Π_λ can be written as $C_\lambda(x) = [L_\lambda(x), R_\lambda(x)]$, with $L_\lambda(x) = (x - \lambda^{-1}E_L) \vee 0$ and $R_\lambda(x) = (x + \lambda^{-1}E_R) \wedge 1$, where E_L, E_R are two independent $\text{Exp}(1)$ random variables. Now, since $X \sim \mathcal{U}([0, 1])$ and $f(u) = 1 + u$,

$$\bar{f}_\lambda(x) = \frac{1}{R_\lambda(x) - L_\lambda(x)} \int_{L_\lambda(x)}^{R_\lambda(x)} (1 + u) du = 1 + \frac{L_\lambda(x) + R_\lambda(x)}{2}.$$

Since $L_\lambda(x)$ and $R_\lambda(x)$ are independent, we have

$$\text{Var}(\bar{f}_\lambda(x)) = \frac{\text{Var}(L_\lambda(x)) + \text{Var}(R_\lambda(x))}{4}.$$

In addition,

$$\text{Var}(R_\lambda(x)) = \text{Var}(x + \lambda^{-1}[E_R \wedge \lambda(1 - x)]) = \lambda^{-2}\text{Var}(E_R \wedge [\lambda(1 - x)])$$

Now, if $E \sim \text{Exp}(1)$ and $a \geq 0$, we have

$$(5.3) \quad \begin{aligned} \mathbb{E}[E \wedge a] &= \int_0^a u e^{-u} du + a \mathbb{P}(E \geq a) = 1 - e^{-a} \\ \mathbb{E}[(E \wedge a)^2] &= \int_0^a u^2 e^{-u} du + a^2 \mathbb{P}(E \geq a) = 2(1 - (a+1)e^{-a}), \end{aligned}$$

so that

$$\text{Var}(E \wedge a) = \mathbb{E}[(E \wedge a)^2] - \mathbb{E}[E \wedge a]^2 = 1 - 2ae^{-a} - e^{-2a}.$$

The formula above gives the variances of $R_\lambda(x)$ and $L_\lambda(x)$ respectively:

$$\begin{aligned} \text{Var}(R_\lambda(x)) &= \lambda^{-2}(1 - 2\lambda(1-x)e^{-\lambda(1-x)} - e^{-2\lambda(1-x)}) \\ \text{Var}(L_\lambda(x)) &= \lambda^{-2}(1 - 2\lambda x e^{-\lambda x} - e^{-2\lambda x}), \end{aligned}$$

and thus

$$(5.4) \quad \text{Var}(\bar{f}_\lambda(x)) = \frac{1}{4\lambda^2}(2 - 2\lambda x e^{-\lambda x} - 2\lambda(1-x)e^{-\lambda(1-x)} - e^{-2\lambda x} - e^{-2\lambda(1-x)}).$$

In addition, the formula (5.3) yields

$$\begin{aligned} \mathbb{E}[R_\lambda(x)] &= x + \lambda^{-1}(1 - e^{-\lambda(1-x)}) \\ \mathbb{E}[L_\lambda(x)] &= x - \lambda^{-1}(1 - e^{-\lambda x}), \end{aligned}$$

and thus

$$(5.5) \quad \tilde{f}_\lambda(x) = 1 + \frac{\mathbb{E}[L_\lambda(x)] + \mathbb{E}[R_\lambda(x)]}{2} = 1 + x + \frac{1}{2\lambda}(e^{-\lambda x} - e^{-\lambda(1-x)}).$$

Combining (5.4) and (5.5) with the decomposition (5.2) gives

$$(5.6) \quad \mathbb{E}[(\bar{f}_\lambda(x) - f(x))^2] = \frac{1}{2\lambda^2} \left(1 - \lambda x e^{-\lambda x} - \lambda(1-x)e^{-\lambda(1-x)} - e^{-\lambda} \right).$$

Integrating over X , we obtain

$$(5.7) \quad \begin{aligned} &\mathbb{E}[(\bar{f}_\lambda(X) - f(X))^2] \\ &= \frac{1}{2\lambda^2} \left(1 - \int_0^1 \lambda x e^{-\lambda x} dx - \int_0^1 \lambda(1-x)e^{-\lambda(1-x)} dx - e^{-\lambda} \right) \\ &= \frac{1}{2\lambda^2} \left(1 - 2 \times \frac{1}{\lambda}(1 - (\lambda+1)e^{-\lambda}) - e^{-\lambda} \right) \\ &= \frac{1}{2\lambda^2} \left(1 - \frac{2}{\lambda} + e^{-\lambda} + \frac{2}{\lambda}e^{-\lambda} \right). \end{aligned}$$

Now, note that the bias $\mathbb{E}[(\bar{f}_\lambda(X) - f(X))^2]$ is positive for $\lambda \in \mathbf{R}_+^*$ (indeed, it is nonnegative, and non-zero since f is not piecewise constant). In addition, the expression (5.7) shows that it is continuous in λ on \mathbf{R}_+^* , and that it admits a limit $\frac{1}{12}$ as $\lambda \rightarrow 0$ (using the fact that $e^{-\lambda} = 1 - \lambda + \frac{\lambda^2}{2} - \frac{\lambda^3}{6} + o(\lambda^3)$). Hence, the function $\lambda \mapsto \mathbb{E}[(\bar{f}_\lambda(X) - f(X))^2]$ is positive and continuous on \mathbf{R}_+ , so that it admits a minimum $C_1 > 0$ on the compact interval $[0, 6]$. In addition, the expression (5.7) shows that for $\lambda \geq 6$, we have

$$(5.8) \quad \mathbb{E}[(\bar{f}_\lambda(X) - f(X))^2] \geq \frac{1}{2\lambda^2} \left(1 - \frac{2}{6}\right) = \frac{1}{3\lambda^2}.$$

First lower bound on the variance. We now turn to the task of bounding the variance from below. In order to avoid restrictive conditions on λ , we will provide two separate lower bounds, valid in two different regimes.

Our first lower bound on the variance, valid for $\lambda \leq n/3$, controls the error of estimation of the optimal labels in nonempty cells. It depends on σ^2 , and is of order $\Theta(\sigma^2 \frac{\lambda}{n})$. We use a general bound on the variance of regressograms [1, Proposition 2] (note that while this result is stated for a fixed number of cells, it can be adapted to a random number of cells by conditioning on $K_\lambda = k$ and then by averaging):

$$(5.9) \quad \begin{aligned} & \mathbb{E}[(\hat{f}_{\lambda,n}(X) - \tilde{f}_\lambda(X))^2] \\ & \geq \frac{\sigma^2}{n} \left(\mathbb{E}[K_\lambda] - 2\mathbb{E}_{\Pi_\lambda} \left[\sum_{\mathbf{v} \in \mathcal{L}(\Pi_\lambda)} \exp(-n\mathbb{P}(X \in C_{\mathbf{v}})) \right] \right). \end{aligned}$$

Now, recall that the splits defining Π_λ form a Poisson point process on $[0, 1]$ of intensity λdx (Fact 1). In particular, the splits can be described as follows. Let $(E_k)_{k \geq 1}$ be an i.i.d. sequence of $\text{Exp}(1)$ random variables, and $S_p := \sum_{k=1}^p E_k$ for $p \geq 0$. Then, the (ordered) splits in Π_λ have the same distribution as $(\lambda^{-1}S_1, \dots, \lambda^{-1}S_{K_\lambda-1})$, where $K_\lambda := 1 + \sup\{p \geq 0 : S_p \leq \lambda\}$. In addition, the probability that $X \sim \mathcal{U}([0, 1])$ falls in the cell

$[\lambda^{-1}S_{k-1}, \lambda^{-1}S_k \wedge 1)$ ($1 \leq k \leq K_\lambda$) is $\lambda^{-1}(S_k \wedge 1 - S_{k-1})$, so that

$$\begin{aligned}
 & \mathbb{E} \left[\sum_{\mathbf{v} \in \mathcal{L}(\Pi_\lambda)} \exp(-n\mathbb{P}(X \in C_{\mathbf{v}})) \right] \\
 &= \mathbb{E} \left[\sum_{k=1}^{K_\lambda-1} e^{-n\lambda^{-1}(S_k - S_{k-1})} + e^{-n(1-\lambda^{-1}S_{K_\lambda-1})} \right] \\
 &\leq \mathbb{E} \left[\sum_{k=1}^{\infty} \mathbf{1}(S_k \leq \lambda) e^{-n\lambda^{-1}E_k} \right] + 1 \\
 (5.10) \quad &= \sum_{k=1}^{\infty} \mathbb{E}[\mathbf{1}(S_k \leq \lambda)] \mathbb{E}[e^{-n\lambda^{-1}E_k}] + 1 \\
 &= \sum_{k=1}^{\infty} \mathbb{E}[\mathbf{1}(S_k \leq \lambda)] \cdot \int_0^\infty e^{-n\lambda^{-1}u} e^{-u} du + 1 \\
 &= \frac{\lambda}{n+\lambda} \mathbb{E} \left[\sum_{k=1}^{\infty} \mathbf{1}(S_k \leq \lambda) \right] + 1 \\
 &= \frac{\lambda}{n+\lambda} \mathbb{E}[K_\lambda] + 1 \\
 (5.11) \quad &= \frac{\lambda}{n+\lambda} (1+\lambda) + 1
 \end{aligned}$$

where (5.10) comes from the fact that E_k and S_{k-1} are independent. Plugging Equation (5.11) in the lower bound (5.9) yields

$$\begin{aligned}
 \mathbb{E} \left[(\widehat{f}_{\lambda,n}(X) - \widetilde{f}_\lambda(X))^2 \right] &\geq \frac{\sigma^2}{n} \left((1+\lambda) - 2(1+\lambda) \frac{\lambda}{n+\lambda} - 2 \right) \\
 &= \frac{\sigma^2}{n} \left((1+\lambda) \frac{n-\lambda}{n+\lambda} - 2 \right).
 \end{aligned}$$

Now, assume that $6 \leq \lambda \leq \frac{n}{3}$. Since

$$(1+\lambda) \frac{n-\lambda}{n+\lambda} - 2 \underset{(\lambda \leq n/3)}{\geq} (1+\lambda) \frac{n-n/3}{n+n/3} - 2 = (1+\lambda) \frac{1}{2} - 2 \underset{(\lambda \geq 6)}{\geq} \frac{\lambda}{4},$$

the above lower bound implies, for $6 \leq \lambda \leq \frac{n}{3}$,

$$(5.12) \quad \mathbb{E} \left[(\widehat{f}_{\lambda,n}(X) - \widetilde{f}_\lambda(X))^2 \right] \geq \frac{\sigma^2 \lambda}{4n}.$$

Second lower bound on the variance. The lower bound (5.12) is only valid for $\lambda \leq n/3$; as λ becomes of order n or larger, the previous bound becomes vacuous. We now provide another lower bound on the variance, valid when $\lambda \geq n/3$, by considering the contribution of empty cells to the variance.

Let $\mathbf{v} \in \mathcal{L}(\Pi_\lambda)$. If $C_{\mathbf{v}}$ contains no sample point from \mathcal{D}_n , then for $x \in C_{\mathbf{v}}$: $\widehat{f}_{\lambda,n}(x) = 0$ and thus $(\widehat{f}_{\lambda,n}(x) - \bar{f}_\lambda(x))^2 = \bar{f}_\lambda(x)^2 \geq 1$. Hence, the variance term is lower bounded as follows, denoting $N_n(C)$ the number of $1 \leq i \leq n$ such that $X_i \in C$ and $N_{\lambda,n}(x) = N_n(C_\lambda(x))$:

$$\begin{aligned}
(5.13) \quad \mathbb{E}[(\widehat{f}_{\lambda,n}(X) - \bar{f}_\lambda(X))^2] &\geq \mathbb{P}(N_{\lambda,n}(X) = 0) \\
&= \mathbb{E}\left[\sum_{\mathbf{v} \in \mathcal{L}(\Pi_\lambda)} \mathbb{P}(X \in C_{\mathbf{v}}) \mathbb{P}(N_n(C_{\mathbf{v}}) = 0)\right] \\
&= \mathbb{E}\left[\sum_{\mathbf{v} \in \mathcal{L}(\Pi_\lambda)} \mathbb{P}(X \in C_{\mathbf{v}}) (1 - \mathbb{P}(X \in C_{\mathbf{v}}))^n\right] \\
&\geq \mathbb{E}\left[\left(\sum_{\mathbf{v} \in \mathcal{L}(\Pi_\lambda)} \mathbb{P}(X \in C_{\mathbf{v}}) (1 - \mathbb{P}(X \in C_{\mathbf{v}}))\right)^n\right]
\end{aligned}$$

$$(5.14) \quad \geq \mathbb{E}\left[\sum_{\mathbf{v} \in \mathcal{L}(\Pi_\lambda)} \mathbb{P}(X \in C_{\mathbf{v}}) (1 - \mathbb{P}(X \in C_{\mathbf{v}}))^n\right]$$

$$(5.15) \quad = \left(1 - \mathbb{E}\left[\sum_{\mathbf{v} \in \mathcal{L}(\Pi_\lambda)} \mathbb{P}(X \in C_{\mathbf{v}})^2\right]\right)^n$$

where (5.13) and (5.14) come from Jensen's inequality applied to the convex function $x \mapsto x^n$. Now, using the notations defined above, we have

$$\begin{aligned}
(5.16) \quad \mathbb{E}\left[\sum_{\mathbf{v} \in \Pi_\lambda} \mathbb{P}(X \in C_{\mathbf{v}})^2\right] &\leq \mathbb{E}\left[\sum_{k=1}^{K_\lambda} (\lambda^{-1} E_k)^2\right] \\
&= \lambda^{-2} \mathbb{E}\left[\sum_{k=1}^{\infty} \mathbf{1}(S_{k-1} \leq \lambda) E_k^2\right] \\
&= \lambda^{-2} \mathbb{E}\left[\sum_{k=1}^{\infty} \mathbf{1}(S_{k-1} \leq \lambda) \mathbb{E}[E_k^2 | S_{k-1}]\right] \\
&= 2\lambda^{-2} \mathbb{E}\left[\sum_{k=1}^{\infty} \mathbf{1}(S_{k-1} \leq \lambda)\right]
\end{aligned}$$

$$\begin{aligned}
(5.17) \quad &= 2\lambda^{-2} \mathbb{E}[K_\lambda] \\
&= \frac{2(\lambda + 1)}{\lambda^2},
\end{aligned}$$

where the equality $\mathbb{E}[E_k^2 | S_{k-1}] = 2$ (used in Equation (5.16)) comes from the fact that $E_k \sim \text{Exp}(1)$ is independent of S_{k-1} .

The bounds (5.15) and (5.17) imply that, if $2(\lambda + 1)/\lambda^2 \leq 1$, then

$$(5.18) \quad \mathbb{E}[(\widehat{f}_{\lambda,n}(X) - \bar{f}_\lambda(X))^2] \geq \left(1 - \frac{2(\lambda + 1)}{\lambda^2}\right)^n.$$

Now, assume that $n \geq 18$ and $\lambda \geq \frac{n}{3} \geq 6$. Then

$$\frac{2(\lambda + 1)}{\lambda^2} \leq 2 \cdot \frac{3}{n} \left(1 + \frac{3}{n}\right) \leq 2 \cdot \frac{3}{n} \left(1 + \frac{3}{18}\right) = \frac{7}{n} \underset{(n \geq 18)}{\leq} 1,$$

so that, using the inequality $(1 - x)^m \geq 1 - mx$ for $m \geq 0$ and $x \in \mathbf{R}$,

$$\left(1 - \frac{2(\lambda + 1)}{\lambda^2}\right)^{n/8} \geq \left(1 - \frac{7}{n}\right)^{n/8} \geq 1 - \frac{n}{8} \cdot \frac{7}{n} = \frac{1}{8}.$$

Combining the above inequality with (5.18) gives, letting $C_2 := 1/8^8$,

$$(5.19) \quad \mathbb{E}[(\widehat{f}_{\lambda,n}(X) - \bar{f}_\lambda(X))^2] \geq C_2.$$

Summing up. Assume that $n \geq 18$. Recall the bias-variance decomposition (5.1) of the risk $R(\widehat{f}_{\lambda,n})$ of the Mondrian tree.

- If $\lambda \leq 6$, we saw that the bias (and hence the risk) is larger than C_1 ;
- If $\lambda \geq \frac{n}{3}$, Equation (5.18) implies that the variance (and hence the risk) is larger than C_2 ;
- If $6 \leq \lambda \leq \frac{n}{3}$, Equations (5.8) (bias term) and (5.12) (variance term) imply that

$$R(\widehat{f}_{\lambda,n}) \geq \frac{1}{3\lambda^2} + \frac{\sigma^2\lambda}{4n}.$$

In particular,

$$(5.20) \quad \inf_{\lambda \in \mathbf{R}^+} R(\widehat{f}_{\lambda,n}) \geq C_1 \wedge C_2 \wedge \inf_{\lambda \in \mathbf{R}^+} \left(\frac{1}{3\lambda^2} + \frac{\sigma^2\lambda}{4n}\right) = C_0 \wedge \frac{1}{4} \left(\frac{3\sigma^2}{n}\right)^{2/3}$$

where we let $C_0 = C_1 \wedge C_2$.

6. Proof of Proposition 4. First, note that in all cases, since $|Y| \leq B$ almost surely, we also have $|\widehat{g}_n(X)| \leq B$ almost surely, so that $(Y - \widehat{g}_n(X))^2 \leq 4B^2$. Let $N_\varepsilon = |I_\varepsilon|$. Note that N_ε is a binomial variable with parameters $n - n_0 \geq n/2$ and $\mathbb{P}(X \in B_\varepsilon) \geq p_0(1 - 2\varepsilon)^d$ (since $p \geq p_0$). Now, recall Chernoff's bound: if $N \sim \text{Bin}(m, p)$ and $\delta \in (0, 1)$, then $\mathbb{P}(N \leq$

$(1 - \delta)mq \leq e^{-mq\delta^2/2}$; in particular, $\mathbb{P}(N \leq mq/2) \leq e^{-mq/8}$. Hence, letting $c_1 = p_0(1 - 2\varepsilon)^d/4$,

$$(6.1) \quad \mathbb{P}(N_\varepsilon \leq c_1 n) \leq \exp(-c_1 n/4).$$

Conditionally on I_ε , the sample $\mathcal{D}' = \{(X_i, Y_i) : i \in I_\varepsilon\}$ is an i.i.d. sample of size N_ε of the conditional distribution of (X, Y) given $X \in B_\varepsilon$; it is also independent of \mathcal{D}_{n_0} , and thus of the estimators \hat{f}_α , $\alpha = 0, \dots, A$. It follows from Theorem 1 in the supplementary material ‘‘Proof of the optimality of the empirical star algorithm’’ of [2] that the estimator \hat{g}_n defined by (5.6) satisfies, with probability $1 - \delta$ over the random sample \mathcal{D}' conditionally on N_ε ,

$$(6.2) \quad \begin{aligned} \mathbb{E}_{(X,Y)}[(\hat{g}_n(X) - Y)^2 | X \in B_\varepsilon] - \min_{0 \leq \alpha \leq A} \mathbb{E}_{(X,Y)}[(\hat{f}_\alpha(X) - Y)^2 | X \in B_\varepsilon] \\ \leq \frac{CB^2 \log[(A+1)\delta^{-1}]}{N_\varepsilon} \end{aligned}$$

for every $\delta \in (0, 1)$, where $C = 600$ and the expectation is taken with respect to an independent sample (X, Y) (the bound (6.2) is deduced from the aforementioned theorem by replacing Y by Y/B , which lies in $[-1, 1]$). Since $Y = f(X) + \varepsilon$ with $\mathbb{E}[\varepsilon | X] = 0$, we have $\mathbb{E}[(g(X) - Y)^2 | X] = \mathbb{E}[(g(X) - f(X))^2 | X] + \mathbb{E}[\varepsilon^2 | X]$. Hence, inequality (6.2) writes

$$\begin{aligned} \mathbb{E}_{(X,Y)}[(\hat{g}_n(X) - f(X))^2 | X \in B_\varepsilon] \\ \leq \min_{0 \leq \alpha \leq A} \mathbb{E}_{(X,Y)}[(\hat{f}_\alpha(X) - f(X))^2 | X \in B_\varepsilon] + \frac{CB^2 \log[(A+1)\delta^{-1}]}{N_\varepsilon}. \end{aligned}$$

By integrating the above inequality over the confidence level δ , we obtain

$$\begin{aligned} \mathbb{E}_{(X,Y), \mathcal{D}'}[(\hat{g}_n(X) - f(X))^2 | X \in B_\varepsilon, N_\varepsilon] \\ \leq \min_{0 \leq \alpha \leq A} \mathbb{E}_{(X,Y)}[(\hat{f}_\alpha(X) - f(X))^2 | X \in B_\varepsilon] + \frac{CB^2[\log(A+1) + 1]}{N_\varepsilon}; \end{aligned}$$

by taking the expectation over \mathcal{D}_{n_0} , conditioning on $N_\varepsilon > c_1 n$, and recalling that $A \leq \log_2(n)$, we get

$$(6.3) \quad \begin{aligned} \mathbb{E}[(\hat{g}_n(X) - f(X))^2 | X \in B_\varepsilon, N_\varepsilon > c_1 n] \\ \leq \min_{0 \leq \alpha \leq A} \mathbb{E}[(\hat{f}_\alpha(X) - f(X))^2 | X \in B_\varepsilon] + \frac{CB^2[\log(1 + \log_2 n) + 1]}{c_1 n}. \end{aligned}$$

Finally, combining the bounds (6.1) and (6.3) yields

$$\begin{aligned}
 & \mathbb{E}[(\widehat{g}_n(X) - f(X))^2 | X \in B_\varepsilon] \\
 & \leq \mathbb{P}(N_\varepsilon \leq c_1 n) \cdot 4B^2 + \mathbb{E}[(\widehat{g}_n(X) - f(X))^2 | X \in B_\varepsilon, N_\varepsilon > c_1 n] \\
 (6.4) \quad & \leq 4B^2 e^{-c_1 n/4} + \min_{0 \leq \alpha \leq A} \mathbb{E}[(\widehat{f}_\alpha(X) - f(X))^2 | X \in B_\varepsilon] \\
 & \quad + \frac{CB^2[\log(1 + \log_2 n) + 1]}{c_1 n},
 \end{aligned}$$

which is precisely inequality (5.7).

Assume that f belongs to the class $\mathcal{C}^{p,\beta}(L)$, with $p \in \{0, 1\}$, $\beta \in (0, 1]$ and $L > 0$; we now proceed to show that \widehat{g}_n achieves the minimax rate of estimation for this class. Let $s = p + \beta \in (0, 2]$. If $p = 0$ (namely, $s \leq 1$), it follows from Theorem 2 (with the same adaptation as in the proof of Theorem 3 to bound the variance term conditionally on $X \in B_\varepsilon$) that, for every $\lambda > 0$,

$$\mathbb{E}[(\widehat{f}_{\lambda, n_0, M}(X) - f(X))^2 | X \in B_\varepsilon] \leq \frac{(4d)^s L^2}{\lambda^{2s}} + \frac{11B^2(1 + \lambda)^d}{p_0(1 - 2\varepsilon)^d n_0}$$

(note that $\sigma, \|f\|_\infty \leq B$ since $|Y| \leq B$). It follows that, for some constants C_1, C_2 independent of λ, L, n ,

$$\begin{aligned}
 \min_{0 \leq \alpha \leq A} \mathbb{E}[(\widehat{f}_\alpha(X) - f(X))^2 | X \in B_\varepsilon] & \leq \min_{0 \leq \alpha \leq A} \left[\frac{C_1 L^2}{(2^\alpha)^{2s}} + \frac{C_2(1 + 2^\alpha)^d}{n} \right] \\
 (6.5) \quad & \leq 4 \min_{\lambda \in [1, n^{1/d}]} \left[\frac{C_1 L^2}{\lambda^{2s}} + \frac{C_2(1 + \lambda)^d}{n} \right],
 \end{aligned}$$

where we used the fact that, for every $\lambda \in [1, n^{1/d}]$, there exists some α , $0 \leq \alpha \leq A$, such that $\lambda/2 \leq 2^\alpha \leq \lambda$. It follows from (6.4) and (6.5) that

$$\begin{aligned}
 & \mathbb{E}[(\widehat{g}_n(X) - f(X))^2 | X \in B_\varepsilon] \\
 & \leq O\left(\min_{0 \leq \lambda \leq n^{1/d}} \left[\frac{C_1 L^2}{\lambda^{2s}} + \frac{C_2(1 + \lambda)^d}{n} \right] + \frac{\log \log n}{n}\right) \\
 & = O\left(L^{2d/(d+2s)} n^{-2s/(d+2s)}\right)
 \end{aligned}$$

where the last bound follows from the fact that $\lambda_* = (L^2 n)^{1/(d+2s)}$ belongs to $[1, n^{1/d}]$ for n large enough (and $\log \log n/n = o(n^{2s/(d+2s)})$).

Now, consider the case $p = 1$, *i.e.*, $1 < s \leq 2$. It follows from Theorem 3 that for some constants C_3, C_4 independent of λ, L, n , we have for every

$\lambda \in [1, n^{1/d}]$ (using the fact that $M \geq n^{2/d} \geq \lambda^2$, so that $1/(M\lambda^2) \leq 1/\lambda^4 \leq 1/\lambda^{2s}$, and $e^{-\lambda\varepsilon}/\lambda^3 = O(1/\lambda^{2s})$)

$$(6.6) \quad \mathbb{E}[(\widehat{f}_{\lambda,n,M}(X) - f(X))^2 | X \in B_\varepsilon] \leq \frac{C_3 L^2}{\lambda^{2s}} + \frac{C_4(1+\lambda)^d}{n}.$$

From the same argument as in the case $0 < s \leq 1$, combining inequalities (6.6) and (6.4) yields

$$\mathbb{E}[(\widehat{g}_n(X) - f(X))^2 | X \in B_\varepsilon] = O(L^{2d/(d+2s)} n^{-2s/(d+2s)})$$

which concludes the proof of Proposition 4. \square

7. Proof of Lemma 1. According to Equation (7.15) from the main text, we have

$$(7.1) \quad F_\lambda(x, z) = \lambda^d \exp(-\lambda\|x - z\|_1) \prod_{1 \leq j \leq d} G_\lambda(x_j, z_j)$$

where we defined, for $u, v \in [0, 1]$,

$$\begin{aligned} G_\lambda(u, v) &= \mathbb{E} \left[(\lambda|u - v| + E_1 \wedge \lambda(u \wedge v) + E_2 \wedge \lambda(1 - u \vee v))^{-1} \right] \\ &= H(\lambda|u - v|, \lambda u \wedge v, \lambda(1 - u \vee v)) \end{aligned}$$

with E_1, E_2 two independent $\text{Exp}(1)$ random variables, and $H : (\mathbf{R}_+^*)^3 \rightarrow \mathbf{R}$ the function defined by

$$H(a, b_1, b_2) = \mathbb{E} \left[(a + E_1 \wedge b_1 + E_2 \wedge b_2)^{-1} \right];$$

also, let

$$H(a) = \mathbb{E} \left[(a + E_1 + E_2)^{-1} \right].$$

Denote

$$\begin{aligned} A &= \int_{[0,1]^d} (z - x) F_\lambda(x, z) dz \\ B &= \int_{[0,1]^d} \frac{1}{2} \|z - x\|^2 F_\lambda(x, z) dz. \end{aligned}$$

Since $1 = \int F_\lambda^{(1)}(u, v) dv = \int \lambda \exp(-\lambda|u - v|) G_\lambda(u, v) dv$, applying Fubini's theorem we obtain

$$(7.2) \quad A_j = \Phi_\lambda^1(x_j) \quad \text{and} \quad B = \sum_{j=1}^d \Phi_\lambda^2(x_j)$$

where we define for $u \in [0, 1]$ and $k \in \mathbf{N}$

$$(7.3) \quad \Phi_\lambda^k(u) = \int_0^1 \lambda \exp(-\lambda|u-v|) G_\lambda(u, v) \frac{(v-u)^k}{k!} dv.$$

Observe that

$$\Phi_\lambda^k(u) = \lambda^{-k} \int_{-\lambda u}^{\lambda(1-u)} \frac{v^k}{k!} \exp(-|v|) H(|v|, \lambda u + v \wedge 0, \lambda(1-u) - v \vee 0) dv.$$

We will control $\Phi_\lambda^k(u)$ for $k = 1, 2$. First, write

$$\begin{aligned} \lambda \Phi_\lambda^1(u) &= - \int_0^{\lambda u} v e^{-v} H(v, \lambda u - v, \lambda(1-u)) dv \\ &\quad + \int_0^{\lambda(1-u)} v e^{-v} H(v, \lambda u, \lambda(1-u) - v) dv \end{aligned}$$

Now, let $\beta := \lambda \frac{u \wedge (1-u)}{2}$. We have

$$\begin{aligned} &\lambda \Phi_\lambda^1(u) - \int_0^\beta v e^{-v} [H(v, \lambda u, \lambda(1-u) - v) - H(v, \lambda u - v, \lambda(1-u))] dv = \\ &- \underbrace{\int_\beta^{\lambda u} v e^{-v} H(v, \lambda u - v, \lambda(1-u)) dv}_{:=I_1 \geq 0} + \underbrace{\int_\beta^{\lambda(1-u)} v e^{-v} H(v, \lambda u, \lambda(1-u) - v) dv}_{:=I_2 \geq 0} \end{aligned}$$

so that the left-hand side of the above equation is between $-I_1 \leq 0$ and $I_2 \geq 0$, and thus its absolute value is bounded by $|I_1| \vee |I_2|$. Now, note that, since $H(v, \cdot, \cdot) \leq v^{-1}$, we have

$$|I_2| \leq \int_\beta^\infty v e^{-v} v^{-1} dv = e^{-\beta}$$

and similarly $|I_1| \leq e^{-\beta}$, so that

$$(7.4) \quad \left| \lambda \Phi_\lambda^1(u) - \underbrace{\int_0^\beta v e^{-v} [H(v, \lambda u, \lambda(1-u) - v) - H(v, \lambda u - v, \lambda(1-u))] dv}_{:=I_3} \right| \leq e^{-\beta}$$

It now remains to bound $|I_3|$. For that purpose, note that since H is decreasing in its second and third argument, we have

$$\begin{aligned} & H(v) - H(v, \lambda u - v, \lambda(1 - u)) \\ & \leq H(v, \lambda u, \lambda(1 - u) - v) - H(v, \lambda u - v, \lambda(1 - u)) \\ & \leq H(v, \lambda u, \lambda(1 - u) - v) - H(v) \end{aligned}$$

which implies

$$\begin{aligned} & |H(v, \lambda u, \lambda(1 - u) - v) - H(v, \lambda u - v, \lambda(1 - u))| \\ & \leq \max(|H(v, \lambda u, \lambda(1 - u) - v) - H(v)|, |H(v) - H(v, \lambda u - v, \lambda(1 - u))|). \end{aligned}$$

Besides, since $(a + E_1 \wedge b_1 + E_2 \wedge b_2)^{-1} \leq (a + E_1 + E_2)^{-1} + a^{-1}(\mathbf{1}\{E_1 \geq b_1\} + \mathbf{1}\{E_2 \geq b_2\})$,

$$(7.5) \quad H(a, b_1, b_2) - H(a) \leq a^{-1}(e^{-b_1} + e^{-b_2}),$$

for all a, b_1, b_2 . Since $\lambda u - v \geq \beta$ and $\lambda(1 - u) - v \geq \beta$ for $v \in [0, \beta]$, we have

$$|H(v) - H(v, \lambda u - v, \lambda(1 - u))|, |H(v) - H(v, \lambda u, \lambda(1 - u) - v)| \leq 2v^{-1}e^{-\beta}$$

so that for $v \in [0, \beta]$

$$|H(v, \lambda u, \lambda(1 - u) - v) - H(v, \lambda u - v, \lambda(1 - u))| \leq 2v^{-1}e^{-\beta}$$

and hence

$$\begin{aligned} |I_3| & \leq \int_0^\beta v e^{-v} |H(v, \lambda u, \lambda(1 - u) - v) - H(v, \lambda u - v, \lambda(1 - u))| dv \\ & \leq \int_0^\beta v e^{-v} 2v^{-1} e^{-\beta} dv \\ & \leq 2e^{-\beta} \int_0^\infty e^{-v} dv \\ (7.6) \quad & = 2e^{-\beta} \end{aligned}$$

Combining Equations (7.4) and (7.6) yields:

$$(7.7) \quad |\Phi_\lambda^1(u)| \leq \frac{3}{\lambda} e^{-\lambda[u \wedge (1-u)]/2}$$

that is,

$$\left\| \int_{[0,1]^d} (z - x) F_\lambda(x, z) dz \right\|^2 = \sum_{j=1}^d (\Phi_\lambda^1(x_j))^2 \leq \frac{9}{\lambda^2} \sum_{j=1}^d e^{-\lambda[x_j \wedge (1-x_j)]}.$$

Furthermore,

$$\begin{aligned} 0 \leq \Phi_\lambda^2(u) &= \lambda^{-2} \int_{-\lambda u}^{\lambda(1-u)} \frac{v^2}{2} e^{-|v|} H(|v|, \lambda u + v \wedge 0, \lambda(1-u) - v \vee 0) dv \\ &\leq \lambda^{-2} \int_0^\infty v^2 e^{-v} v^{-1} dv \\ &= \lambda^{-2} \end{aligned}$$

so that

$$0 \leq \Phi_\lambda^2(u) \leq \frac{1}{\lambda^2},$$

which proves the second inequality by summing over $j = 1, \dots, d$. This concludes the proof of Lemma 1. \square

References.

- [1] S. Arlot and R. Genuer. Analysis of purely random forests bias. *arXiv preprint arXiv:1407.3939*, 2014.
- [2] J.-Y. Audibert. Progressive mixture rules are deviation suboptimal. In *Advances in Neural Information Processing Systems 20*, pages 41–48. Curran Associates, Inc., 2008.
- [3] L. Györfi, M. Kohler, A. Krzyzak, and H. Walk. *A distribution-free theory of nonparametric regression*. Springer Science & Business Media, 2002.

JAOUAD MOURTADA
CMAP, ECOLE POLYTECHNIQUE
ROUTE DE SACLAY
91128 PALAISEAU CEDEX, FRANCE
E-MAIL: jaouad.mourtada@polytechnique.edu

STÉPHANE GAÏFFAS
LPMA - UNIV. PARIS DIDEROT
BÂTIMENT SOPHIE GERMAIN
CASE COURRIER 7012
75205 PARIS CEDEX 13, FRANCE
E-MAIL: stephane.gaiffas@lpsm.paris

ERWAN SCORNET
CMAP, ECOLE POLYTECHNIQUE
ROUTE DE SACLAY
91128 PALAISEAU CEDEX, FRANCE
E-MAIL: erwan.scornet@polytechnique.edu