

---

# Provable Guarantees for Decision Tree Induction: The Agnostic Setting

---

Guy Blanc<sup>\*1</sup> Jane Lange<sup>\*1</sup> Li-Yang Tan<sup>\*1</sup>

## Abstract

We give strengthened provable guarantees on the performance of widely employed and empirically successful *top-down decision tree learning heuristics*. While prior works have focused on the realizable setting, we consider the more realistic and challenging *agnostic* setting. We show that for all monotone functions  $f$  and  $s \in \mathbb{N}$ , these heuristics construct a decision tree of size  $s^{\tilde{O}((\log s)/\varepsilon^2)}$  that achieves error  $\leq \text{opt}_s + \varepsilon$ , where  $\text{opt}_s$  denotes the error of the optimal size- $s$  decision tree for  $f$ . Previously such a guarantee was not known to be achievable by any algorithm, even one that is not based on top-down heuristics. We complement our algorithmic guarantee with a near-matching  $s^{\tilde{\Omega}(\log s)}$  lower bound.

## 1. Introduction

This paper is motivated by the goal of establishing strong provable guarantees for the class of popular and empirically successful *top-down decision tree learning heuristics*. This includes well-known instantiations such as ID3 (Quinlan, 1986), its successor C4.5 (Quinlan, 1993), and CART (Breiman, 2017), all widely employed in everyday machine learning applications. These simple heuristics are also at the heart of more sophisticated decision-tree-based algorithms such as random forests (Breiman, 2001) and XGBoost (Chen & Guestrin, 2016), which have quickly gained prominence in Kaggle and other data science competitions, and achieve state-of-the-art performance for diverse tasks.

We will soon formally describe the learning-theoretic framework within which we study these heuristics, mentioning for

---

<sup>\*</sup>Equal contribution <sup>1</sup>Stanford University. Correspondence to: Guy Blanc <gblanc@stanford.edu>, Jane Lange <jlange20@stanford.edu>, Li-Yang Tan <liyong@cs.stanford.edu>.

now that they build a decision tree  $T$  for a binary classifier  $f : \mathbb{R}^n \rightarrow \{0, 1\}$  in a *greedy, top-down* fashion:

1. Query  $\mathbb{1}[x_i \geq \theta]$  at the root of  $T$ , where  $x_i$  and  $\theta$  are chosen to maximize the *purity gain*

$$\mathcal{G}(\mathbb{E}[f]) - \left( \Pr[x_i \geq \theta] \cdot \mathcal{G}(\mathbb{E}[f_{x_i \geq \theta}]) + \Pr[x_i < \theta] \cdot \mathcal{G}(\mathbb{E}[f_{x_i < \theta}]) \right)$$

with respect to an *impurity function*  $\mathcal{G} : [0, 1] \rightarrow [0, 1]$ . This carefully chosen function  $\mathcal{G}$  encapsulates the *splitting criterion* of the heuristic.

2. Build the left and right subtrees of  $T$  by recursing on  $f_{x_i \geq \theta}$  and  $f_{x_i < \theta}$  respectively.

Different instantiations of this simple approach are distinguished by different impurity functions  $\mathcal{G}$ , which determine the *order* in which the recursive calls are made. For example, ID3 and C4.5 uses the binary entropy function  $\mathcal{G}(p) = H(p)$ ; CART uses the *Gini criterion*  $\mathcal{G}(p) = 4p(1-p)$ ; Kearns and Mansour proposed and analyzed the function  $\mathcal{G}(p) = 2\sqrt{p(1-p)}$  (Kearns & Mansour, 1999; Dietterich et al., 1996).

Even without specifying the impurity function  $\mathcal{G}$ , it is well known and easy to see that *any* such heuristic can fare poorly even for simple functions  $f$ , in the sense of building a decision tree that is much larger than the optimal one. Consider the most basic setting of binary features, the uniform distribution over inputs, and  $f$  being the parity of two variables  $f(x) = x_1 \oplus x_2$ . This function can be computed by a decision tree of size 4, but since  $\mathbb{E}[f] = \mathbb{E}[f_{x_i \geq \theta}] = \mathbb{E}[f_{x_i < \theta}] = 0$  for all  $x_i$  and  $\theta$ , *any* top-down heuristic—regardless of the impurity function  $\mathcal{G}$ —may build a tree of size  $\Omega(2^n)$  before achieving any non-trivial accuracy.

**Monotonicity and the conjectures of Fiat, Pechyony, and Lee.** In light of such examples, a question suggests itself: can we identify natural and expressive classes of functions for which strong provable guarantees on the performance of these top-down heuristics *can* be obtained?

The first results in this direction were given by Fiat and Pechyony (Fiat & Pechyony, 2004; Pechyony, 2004), who considered the case of binary features (i.e.  $f : \{\pm 1\}^n \rightarrow \{0, 1\}$ ) and showed a strong positive result for halfspaces and read-once DNF formulas. For such functions  $f$ , they showed that these heuristics build a decision tree of *optimal* size that compute  $f$  *exactly*. Furthermore, and most relevant to our work, they raised the intriguing possibility that *monotonicity* is the key property shared by these functions that enables such a guarantee.<sup>1</sup> They conjectured that for all monotone functions  $f$ , these heuristics build a decision tree of size “not far from minimal” that compute  $f$  exactly.

Subsequently, Lee (Lee, 2009) formulated a relaxation of Fiat and Pechyony’s conjecture, allowing for an approximate representation rather than an exact representation. Lee conjectured that for all monotone functions  $f : \{\pm 1\}^n \rightarrow \{0, 1\}$  computable by size- $s$  decision trees, these heuristics construct a decision tree of size  $\text{poly}(s, 1/\varepsilon)$  that achieves error  $\varepsilon$  with respect to the uniform distribution over inputs. (The author further remarked that even a bound that is  $\text{poly}(s)$  for constant values of  $\varepsilon$  “would be a huge advance.”)

These conjectures of (Fiat & Pechyony, 2004) and (Lee, 2009) are especially appealing because monotonicity, beyond just being a natural assumption that excludes parity and “parity-like” functions, is an independently important and intensively-studied property in both the theory and practice of machine learning. Many real-world data sets are naturally monotone in their features. In learning theory, even restricting our attention just to uniform-distribution learning, there is a large body of work on learning monotone functions (Hancock & Mansour, 1991; Kearns & Valiant, 1994; Kearns et al., 1994; Bshouty, 1995; Bshouty & Tamon, 1996; Blum et al., 1998; Verbeurgt, 1998; Sakai & Maruoka, 2000; Servedio, 2004; O’Donnell & Servedio, 2007; Sellie, 2008; Dachman-Soled et al., 2009; Lee, 2009; Jackson et al., 2011; O’Donnell & Wimmer, 2013; Dachman-Soled et al., 2015). Partly responsible for this popularity is the fact that monotonicity allows one to sidestep the well-known *statistical query* lower bounds (Blum et al., 1994) that hold for many simple concept classes.

**The work of (Blanc et al., 2020).** Recent work of Blanc, Lange, and Tan established a weak version of (Lee, 2009)’s conjecture. For all monotone functions  $f : \{\pm 1\}^n \rightarrow \{0, 1\}$  that are computable by size- $s$  decision trees, they showed that a close variant of these top-down heuristics constructs a decision tree of size  $s^{O(\sqrt{\log s}/\varepsilon)}$  that achieves error  $\varepsilon$  with respect to the uniform distribution over inputs.

<sup>1</sup>We consider a function to be monotone if it is either non-decreasing or non-increasing in every coordinate; we do not require the direction to be the same for all coordinates. (Such functions are sometimes also called “unate.”) See Definition 5.

(Blanc et al., 2020) also showed that the dependence on ‘ $s$ ’ *cannot* be made polynomial, thereby disproving (Lee, 2009)’s actual conjecture (and (Fiat & Pechyony, 2004)’s even stronger conjecture): for all sizes  $s \leq 2^{\tilde{O}(n^{4/5})}$  and error parameters  $\varepsilon \in (0, \frac{1}{2})$ , they exhibited a monotone function  $f$  that is computable by a size- $s$  decision tree, and showed that all top-down impurity-based heuristics have to build a decision tree of size  $s^{\tilde{\Omega}(\sqrt[4]{\log s})}$  in order to achieve error  $\varepsilon$ .

**Our contributions.** We give strengthened provable guarantees on the performance of top-down decision tree learning heuristics. The three main contributions of our work are:

1. We consider the more realistic and challenging *agnostic* setting, where  $f$  is an *arbitrary* monotone function. Prior works focused on the realizable setting, and their results relied on assumptions about the computational complexity of  $f$  (i.e. that  $f$  is computable by a small decision tree, or a halfspace, or a read-once DNF formula, etc.).
2. We establish provable guarantees that apply to all top-down heuristics, including ID3, C4.5, and CART. (Blanc et al., 2020)’s analysis, on the other hand, dealt with a specific *variant* of these heuristics, one whose splitting criterion does not correspond to any impurity function  $\mathcal{G}$ . (As a secondary contribution, we further show that (Blanc et al., 2020)’s guarantees in the realizable setting also hold for all top-down heuristics.)
3. Our analysis extends to classifiers for *real-valued* features (i.e.  $f : \mathbb{R}^n \rightarrow \{0, 1\}$ ) and arbitrary product distributions over inputs, whereas prior works dealt with classifiers for binary features (i.e.  $f : \{\pm 1\}^n \rightarrow \{0, 1\}$ ) and mostly focused on the uniform distribution over inputs. Trees for real-valued features branch on queries of the form  $\mathbb{1}[x_i \geq \theta]$  for some  $\theta \in \mathbb{R}$ , whereas trees for binary features branch on queries of the form  $\mathbb{1}[x_i = 1]$ .

Our main result is as follows:

**Theorem 1** (Our main result; informal). *Let  $f : \mathbb{R}^n \rightarrow \{0, 1\}$  be a monotone function and  $\mathcal{D}$  be a product distribution over  $\mathbb{R}^n$ . For  $s \in \mathbb{N}$ , let  $\text{opt}_s$  denote the error of the best size- $s$  balanced decision tree for  $f$  with respect to  $\mathcal{D}$ .<sup>2</sup> Let  $\mathcal{G}$  be any impurity function. For  $t \leq s^{\tilde{O}((\log s)/\varepsilon^2)}$ ,*

<sup>2</sup>A balanced decision tree of size  $s$  is one that has depth  $O(\log s)$ . This technical assumption is not necessary in the case of binary features and the uniform distribution over inputs (see Theorem 2); it is only necessary for the general setting of real-valued features and arbitrary product distributions.

the size- $t$  decision tree for  $f$  constructed by the top-down heuristic with  $\mathcal{G}$  as its splitting criterion achieves error  $\leq \text{opt}_s + \varepsilon$  with respect to  $\mathcal{D}$ .

Previously, such a guarantee was not known to be achievable by any algorithm, even one that is not based on top-down heuristics (i.e. [Theorem 1](#) represents the first algorithm for properly learning decision trees in an agnostic setting).

We complement [Theorem 1](#) with a near-matching lower bound. We show that for all  $s \leq 2^{\tilde{O}(\sqrt{n})}$ , there is a monotone function  $f$  such that  $\text{opt}_s \leq 0.01$ , and yet any top-down heuristic has to grow a tree of size  $s^{\tilde{\Omega}(\log s)}$  to even achieve error  $\leq 0.49$  with respect to  $f$ . Taken together with ([Blanc et al., 2020](#))’s  $s^{O(\sqrt{\log s})}$  upper bound in the realizable setting, this exhibits a separation between the realizable and agnostic settings.

### 1.1. Formal Statements of our Results

We define a *partial tree* to be a decision tree with unlabeled leaves, and write  $T^\circ$  to denote such trees. We refer to any decision tree  $T$  obtained from  $T^\circ$  by a labeling of its leaves as a *completion* of  $T^\circ$ . Given a tree a partial tree  $T^\circ$  and a function  $f$ , there is a canonical completion of  $T^\circ$  that minimizes the approximation error with respect to  $f$ :

**Definition 1** ( $f$ -completion of a partial tree). *Let  $T^\circ$  be a partial tree and  $f : \mathbb{R}^n \rightarrow \{0, 1\}$ . Consider the following completion of  $T^\circ$ : for every leaf  $\ell$  in  $T^\circ$ , label it  $\text{round}(\mathbb{E}[f_\ell])$ , where  $f_\ell$  denotes the restriction of  $f$  by the path leading to  $\ell$  and  $\text{round}(p) = 1$  if  $p \geq \frac{1}{2}$  and 0 otherwise. This completion minimizes the approximation error  $\Pr[T(\mathbf{x}) \neq f(\mathbf{x})]$ ; we refer to it as the  $f$ -completion of  $T^\circ$  and denote it as  $T_f^\circ$ .*

**Definition 2** (Impurity functions and strong concavity). *An impurity function  $\mathcal{G} : [0, 1] \rightarrow [0, 1]$  is a concave function that is symmetric around  $\frac{1}{2}$ , and satisfies  $\mathcal{G}(0) = \mathcal{G}(1) = 0$  and  $\mathcal{G}(\frac{1}{2}) = 1$ . We say that  $\mathcal{G}$  is  $\kappa$ -strongly concave if for all  $a, b \in [0, 1]$ ,*

$$\frac{\mathcal{G}(a) + \mathcal{G}(b)}{2} \leq \mathcal{G}\left(\frac{a+b}{2}\right) - \frac{\kappa}{2} \cdot (b-a)^2.$$

**Remark 1** ( $\kappa$  values of common impurity functions). ID3 and C4.5 use binary entropy as their impurity function, which is  $\kappa$ -strongly concave for  $\kappa = 1/\ln(2)$ , or  $\approx 1.4$ . Gini impurity, which is used by CART, is strongly concave for  $\kappa = 2$ , and ([Kearns & Mansour, 1999](#))’s impurity function is strongly concave for  $\kappa = 1$ .

We now formally define the top-down heuristics that we study.

**Definition 3** ( $\mathcal{G}$ -impurity of a partial tree). *Let  $\mathcal{G} : [0, 1] \rightarrow [0, 1]$  be an impurity function and  $\mathcal{D}$  be a distribution over*

$\mathbb{R}^n$ . For  $f : \mathbb{R}^n \rightarrow \{0, 1\}$  and a partial tree  $T^\circ$ , the  $\mathcal{G}$ -impurity of  $T^\circ$  with respect to  $f$  is defined to be

$$\begin{aligned} \mathcal{G}\text{-impurity}_{f, \mathcal{D}}(T^\circ) \\ := \sum_{\text{leaves } \ell \in T^\circ} \Pr_{\mathbf{x} \sim \mathcal{D}}[\mathbf{x} \text{ reaches } \ell] \cdot \mathcal{G}(\mathbb{E}[f_\ell]). \end{aligned}$$

If  $T^\circ$  is a partial tree and  $\ell$  is a leaf of  $T^\circ$ , we write  $T_{\ell, \mathbb{1}[x_i \geq \theta]}^\circ$  to denote the extension of  $T^\circ$  obtained by splitting  $\ell$  with a query to  $\mathbb{1}[x_i \geq \theta]$ . The following algorithm captures the top-down decision tree learning heuristics that we study in this work:

**BUILDTOPDOWNDT** $_{\mathcal{G}, \mathcal{D}}(f, t)$ :

Initialize  $T^\circ$  to be the empty tree.

while ( $\text{size}(T^\circ) < t$ ) {

    Grow  $T^\circ$  by splitting leaf  $\ell$  with a query to  $\mathbb{1}[x_i \geq \theta]$ , where  $\ell$  and  $\mathbb{1}[x_i \geq \theta]$  maximize:

$$\begin{aligned} &\mathcal{G}\text{-impurity}_{f, \mathcal{D}}(T^\circ) \\ &- \mathcal{G}\text{-impurity}_{f, \mathcal{D}}(T_{\ell, \mathbb{1}[x_i \geq \theta]}^\circ), \end{aligned}$$

    the *purity gain* with respect to  $\mathcal{G}$  and  $\mathcal{D}$ .

}

Output the  $f$ -completion of  $T^\circ$ .

Figure 1. Top-down heuristic for building a size- $t$  decision tree approximation of a function  $f : \mathbb{R}^n \rightarrow \{0, 1\}$ , using the impurity function  $\mathcal{G}$  as its splitting criterion.

We write  $\text{TOPDOWNERROR}_{\mathcal{G}, \mathcal{D}}(f, t)$  to denote the error  $\Pr[T(\mathbf{x}) \neq f(\mathbf{x})]$ , where  $T$  is the size- $t$  tree constructed by  $\text{BUILDTOPDOWNDT}_{\mathcal{G}, \mathcal{D}}(f, t)$ .

**Definition 4** ( $\text{opt}_s$ ). *For a function  $f : \mathbb{R}^n \rightarrow \{0, 1\}$ , a distribution  $\mathcal{D}$  over  $\mathbb{R}^n$ , and an integer  $s \in \mathbb{N}$ , we write  $\text{opt}_{f, \mathcal{D}, s} \in [0, \frac{1}{2}]$  to denote the error of the best size- $s$  decision tree for  $f$ :*

$\text{opt}_{f, \mathcal{D}, s}$

$$:= \min \left\{ \Pr_{\mathbf{x} \sim \mathcal{D}}[T(\mathbf{x}) \neq f(\mathbf{x})] : T \text{ is a size-}s \text{ decision tree} \right\}.$$

When  $f$  and  $\mathcal{D}$  are clear from context, we simply write  $\text{opt}_s$ .

**Our main algorithmic guarantee.** We first state and prove our results in the setting of binary features and the uniform distribution over inputs. As alluded to in the introduction, in [Section 2.2](#) we will show how our results in this specialized setting extend to the more general setting of real-valued features and arbitrary product distributions over inputs.

**Theorem 2** (Our main algorithmic guarantee for binary features and the uniform distribution). *Let  $f : \{\pm 1\}^n \rightarrow \{0, 1\}$  be a monotone function,  $\mathcal{U}$  be the uniform distribution over  $\{\pm 1\}^n$ , and  $\mathcal{G} : [0, 1] \rightarrow [0, 1]$  be an  $\kappa$ -strongly concave impurity function. For all  $s \in \mathbb{N}$  and  $\varepsilon \in (0, \frac{1}{2})$ ,*

$$\text{TOPDOWNERROR}_{\mathcal{G}, \mathcal{U}}(f, s^{O(\log s)/\kappa\varepsilon^2}) \leq \text{opt}_s + \varepsilon.$$

In words, [Theorem 2](#) says that for  $s \in \mathbb{N}$ , the decision tree of size  $s^{O(\log s)}$  constructed by `BUILDTOPDOWNDT` achieves error that nearly matches that of the best size- $s$  decision tree for  $f$ .

**A near-matching lower bound.** We contrast [Theorem 2](#) with ([Blanc et al., 2020](#))’s result in the *realizable* setting: if  $f$  is a monotone function that is computable by a size- $s$  decision tree (i.e.  $\text{opt}_s = 0$ ), then a variant of the top-down heuristics grows a tree of size  $s^{O(\sqrt{\log s})}$  that achieves error  $\varepsilon$ .<sup>3</sup> With this in mind, it is natural to wonder if the parameters of [Theorem 2](#) can be improved to  $\text{TOPDOWNERROR}_{\mathcal{G}, \mathcal{U}}(f, s^{O_{\kappa, \varepsilon}(\sqrt{\log s})}) \leq \text{opt}_s + \varepsilon$ . We complement [Theorem 2](#) with a lower bound that rules out such an improvement. We show that the dependence on ‘ $s$ ’ in [Theorem 2](#) is in fact near-optimal:

**Theorem 3** (Our main lower bound). *For all  $s \leq 2^{\tilde{O}(\sqrt{n})}$ , there is a monotone function  $f : \{\pm 1\}^n \rightarrow \{0, 1\}$  such that  $\text{opt}_s \leq 0.01$  with respect to the uniform distribution over inputs, and  $\text{TOPDOWNERROR}_{\mathcal{G}, \mathcal{U}}(f, s^{\tilde{\Omega}(\log s)}) \geq 0.49$  for any impurity function  $\mathcal{G}$ .*

Taken together with ([Blanc et al., 2020](#))’s  $s^{O(\sqrt{\log s})}$  upper bound in the realizable setting, [Theorem 3](#) exhibits a separation between the realizable and agnostic settings.

## 1.2. Related Work

Kearns and Mansour ([Kearns, 1996](#); [Kearns & Mansour, 1999](#)) were the first to study top-down decision tree learning heuristics using the framework of learning theory. They showed that these heuristics can be viewed as *boosting algorithms*, where one views the functions queried at the internal nodes of the tree (single variables in our case) as weak hypotheses. As is standard in results on boosting, their results are *conditional* in nature: they assume the existence of weak hypotheses for filtered-and-rebalanced versions of the original distribution (what they call “The Weak Hypothesis Assumption”), and they show how these top-down heuristics build a decision tree that combines these weak hypotheses

<sup>3</sup>The heuristic that ([Blanc et al., 2020](#)) analyzes does not correspond to any impurity function  $\mathcal{G}$ . As mentioned in the introduction, in [Appendix C](#) we show that ([Blanc et al., 2020](#))’s guarantees on their variant of the top-down heuristics in fact hold for all top-down heuristics.

into a strong one. Dietterich, Kearns, and Mansour ([Dietterich et al., 1996](#)) gave an experimental comparison of the impurity functions used by ID3, C4.5, and CART, along with a new impurity function that ([Kearns & Mansour, 1999](#)) had proposed.

Fiat and Pechyony ([Fiat & Pechyony, 2004](#); [Pechyony, 2004](#)) studied functions  $f$  that are computable by linear threshold functions or read-once DNF formulas, and showed that these heuristics build a decision tree of optimal size that compute  $f$  exactly. Recent work of Brutzkus, Daniely, and Malach ([Brutzkus et al., 2019a](#)) studies functions  $f$  that are conjunctions and read-once DNF formulas, and gives theoretical and empirical evidence showing that for such functions, the variant of ID3 proposed by ([Kearns & Mansour, 1999](#)), when run for  $t$  iterations, grows a tree that achieves accuracy that matches or nearly matches that of the best size- $t$  tree for  $f$ . Concurrent work by the same authors ([Brutzkus et al., 2019b](#)) shows that ID3 learns  $(\log n)$ -juntas in the setting of smoothed analysis.

## Learning algorithms not based on top-down heuristics: improper algorithms for learning decision trees.

O’Donnell and Servedio ([O’Donnell & Servedio, 2007](#)) gave a poly( $n, s^{1/\varepsilon^2}$ )-time uniform-distribution algorithm for learning monotone functions computable by size- $s$  decision trees. This remains the fastest algorithm for the realizable setting. ([O’Donnell & Servedio, 2007](#))’s algorithm is not based on the top-down heuristics that are the focus of our work (and the others discussed above); indeed, their algorithm does not output a decision tree as its hypothesis (i.e. it is not a proper learning algorithm). For the agnostic setting, the results of Kalai, Klivans, Mansour, and Servedio ([Kalai et al., 2008](#)) can be used to give a uniform-distribution algorithm that runs in time  $n^{O(\log(s/\varepsilon))}$  and outputs a hypothesis that achieves error  $\text{opt}_s + \varepsilon$ . Compared to [Theorem 2](#), this algorithm does not require  $f$  to be monotone, but like ([O’Donnell & Servedio, 2007](#))’s algorithm it is also improper. The work of ([Gopalan et al., 2008](#)) gives a uniform-distribution algorithm that runs in poly( $n, s, 1/\varepsilon$ ) time; however, their algorithm requires the use of *membership queries*, and is also improper. Furthermore, all the results discussed in this paragraph only hold in the setting of binary features and with respect to the uniform distribution over inputs.

## 1.3. Preliminaries

We use **boldface** (e.g.  $\mathbf{x} \sim \mathbb{R}^n$ ) to denote random variables. Given two functions  $f, g : \mathbb{R}^n \rightarrow \{0, 1\}$ , we write  $\text{dist}(f, g) := \Pr[f(\mathbf{x}) \neq g(\mathbf{x})]$  to denote the distance between  $f$  and  $g$ . We write  $\text{bias}(f) := \min\{\Pr[f(\mathbf{x}) = 0], \Pr[f(\mathbf{x}) = 1]\}$  to denote the distance of  $f$  to the closest constant function; equivalently,  $\text{bias}(f) = \Pr[f(\mathbf{x}) \neq \text{round}(\mathbb{E}[f])]$ , where  $\text{round}(p) = 1$  if  $p \geq \frac{1}{2}$  and 0 other-

wise. If  $\ell$  is a leaf in a decision tree, we write  $|\ell|$  to denote the depth of  $\ell$  within the tree.

**Definition 5** (Monotone functions). *We say that a function  $f : \mathbb{R}^n \rightarrow \{0, 1\}$  is monotone if for all coordinates  $i \in [n]$ , either*

- *$f$  is non-decreasing in the  $i$ -th direction:  $f(x) \leq f(y)$  for all  $x, y \in \mathbb{R}^n$  such that  $x_i \leq y_i$ , or*
- *$f$  is non-increasing in the  $i$ -th direction:  $f(x) \geq f(y)$  for all  $x, y \in \mathbb{R}^n$  such that  $x_i \leq y_i$ .*

**Organization of this paper.** We give the complete proof of [Theorem 2](#) in the next section. Due to space constraints, proofs of the extension of [Theorem 2](#) to [Theorem 1](#) (from the specific setting of binary features and the uniform distribution over inputs, to the general setting of real-valued features and arbitrary product distributions over inputs), and that of [Theorem 3](#) (a lower bound showing that [Theorem 1](#) is nearly optimal) are deferred to the appendix.

## 2. Our Main Algorithmic Result: [Theorem 2](#)

Recall that [Theorem 2](#) is concerned with the special case of binary features and the uniform distribution over inputs. Therefore the trees that we reason about in the proof of [Theorem 2](#) split on queries of the form  $\mathbb{1}[x_i = 1]$  (rather than queries of the form  $\mathbb{1}[x_i \geq \theta]$  as in the general setting of real-valued features). Also, all probabilities and expectations in this proof are with respect to the uniform distribution over inputs.

**Definition 6** (Influence). *Given a function  $f : \{\pm 1\}^n \rightarrow \{0, 1\}$ , the influence of coordinate  $i \in [n]$  on  $f$  is defined to be*

$$\text{Inf}_i(f) := \Pr[f(\mathbf{x}) \neq f(\mathbf{x}^{\oplus i})],$$

where  $\mathbf{x}^{\oplus i}$  denotes  $\mathbf{x}$  with its  $i$ -coordinate flipped, and  $\mathbf{x} \sim \{\pm 1\}$  is uniform random. The total influence of  $f$  is defined to be  $\text{Inf}(f) := \sum_{i=1}^n \text{Inf}_i(f)$ .

A key technical ingredient in our proof will be an inequality of Jain and Zhang ([Jain & Zhang, 2011](#)), a robust version of the powerful O’Donnell–Saks–Schramm–Servedio inequality from the Fourier analysis of boolean functions ([O’Donnell et al., 2005](#)). The following is a special case of Corollary 1.4 of ([Jain & Zhang, 2011](#)):

**Theorem 4** (Robust OSSS inequality). *Let  $f : \{\pm 1\}^n \rightarrow \{0, 1\}$  be any function, and  $g : \{\pm 1\}^n \rightarrow \{0, 1\}$  be a size- $s$  decision tree. Then*

$$\max_{i \in [n]} \{\text{Inf}_i(f)\} \geq \frac{\text{bias}(f) - \text{dist}(f, g)}{\log s}.$$

**Remark 2** (Context for [Theorem 4](#)). The original OSSS inequality essentially corresponds to the special case of [Theorem 4](#) where  $f \equiv g$ : if  $f$  is a size- $s$  decision tree, then

$$\max_{i \in [n]} \{\text{Inf}_i(f)\} \geq \frac{\text{Var}(f)}{\log s}.$$

The OSSS inequality can be viewed as a variant of the famed Kahn–Kalai–Linial inequality ([Kahn et al., 1988](#)), one that takes into account the computational complexity of  $f$ . In ([O’Donnell et al., 2005](#)) the authors also gave a robust version of their inequality that is qualitative similar to [Theorem 4](#) (see the discussion following [Theorem 3.2](#) of ([O’Donnell et al., 2005](#))): under the assumptions of [Theorem 4](#),

$$\max_{i \in [n]} \{\text{Inf}_i(f)\} \geq \frac{\text{Var}(f) - 2 \cdot \text{dist}(f, g)}{\log s}.$$

This inequality can be used in place of [Theorem 4](#) to prove a statement that is qualitatively similar to [Theorem 2](#), but with a weaker error bound of  $O(\text{opt}_s) + \varepsilon$  rather than the  $\text{opt}_s + \varepsilon$  of [Theorem 2](#).

**Fact 2.1** (Influence  $\equiv$  correlation for monotone functions). *Let  $f : \{\pm 1\}^n \rightarrow \{0, 1\}$  be a monotone function. Then  $\text{Inf}_i(f) = 2 \cdot |\mathbb{E}[f(\mathbf{x})\mathbf{x}_i]|$  for all  $i \in [n]$ .*

**Proposition 2.2** (Influence and purity gain). *For all  $\kappa$ -strongly concave impurity heuristics  $\mathcal{G}$ , all monotone functions  $f : \{\pm 1\}^n \rightarrow \{0, 1\}$ , and all coordinates  $i \in [n]$ ,*

$$\mathbb{E}_{\mathbf{b} \sim \{\pm 1\}} [\mathcal{G}(\mathbb{E}[f_{x_i=\mathbf{b}}])] \leq \mathcal{G}(\mathbb{E}[f]) - \frac{\kappa}{32} \cdot \text{Inf}_i(f)^2.$$

*Proof.* Note that:

$$\mathbb{E}[f(\mathbf{x})\mathbf{x}_i] = \frac{1}{2}(\mathbb{E}[f_{x_i=1}] - \mathbb{E}[f_{x_i=-1}])$$

and that:

$$\mathbb{E}[f] = \frac{1}{2}(\mathbb{E}[f_{x_i=1}] + \mathbb{E}[f_{x_i=-1}]).$$

The desired result therefore holds as a direct consequence of the  $\kappa$ -strong concavity of  $\mathcal{G}$  and [Fact 2.1](#).  $\square$

### 2.1. Proof of [Theorem 2](#)

Let  $f : \{\pm 1\}^n \rightarrow \{0, 1\}$  be a monotone function and  $g$  be a size- $s$  decision tree for which  $\text{dist}(f, g) = \text{opt}_s$ . Fix some  $\kappa$ -strongly concave impurity heuristic  $\mathcal{G}$ . Given a partial tree  $T^\circ$  (which we should think of as the approximator for  $f$  that is being built by BUILDTOPDOWNDT $_{\mathcal{G}, \mathcal{U}}$ ), we define the potential function:

$$\mathcal{G}\text{-impurity}_f(T^\circ) := \sum_{\text{leaves } \ell \in T^\circ} \left( 2^{-|\ell|} \cdot \mathcal{G}(\mathbb{E}[f_\ell]) \right).$$

The next lemma records a few useful properties of this potential function  $\mathcal{G}\text{-impurity}_f$ . (This lemma can be viewed

as our analogue of (Blanc et al., 2020)’s Lemma 5.1. (Blanc et al., 2020) worked with a potential function that is a variant of ours which has  $\text{Inf}(f_\ell)$  in place of  $\mathcal{G}(\mathbb{E}[f_\ell])$ .)

**Definition 7** (Purity gain). *Let  $f : \{\pm 1\}^n \rightarrow \{0, 1\}$  be a function and  $\mathcal{G}$  be an impurity heuristic. For a partial tree  $T^\circ$  and a leaf  $\ell$  of  $T^\circ$ , we write  $T_{\ell, x_i}^\circ$  to denote the extension of  $T^\circ$  obtained by splitting  $\ell$  with a query to  $x_i$ . The purity gain associated with splitting  $\ell$  with a query to  $x_i$  is defined as*

$$\begin{aligned} & \mathcal{G}\text{-purity-gain}_f(T^\circ, \ell, x_i) \\ & := \mathcal{G}\text{-impurity}_f(T^\circ) - \mathcal{G}\text{-impurity}_f(T_{\ell, x_i}^\circ). \end{aligned}$$

It is easy to verify that given a leaf  $\ell$ , the variable associated with the largest purity gain is exactly one that is most correlated with  $f_\ell$ :

**Proposition 2.3** (Proposition 7.7 from (Blanc et al., 2020)). *For any leaf  $\ell$  of  $T^\circ$ , let  $x_{i^*}$  be the variable that maximizes  $\mathcal{G}\text{-purity-gain}_{\mathcal{G}}(\ell, x_i)$  among all  $i \in [n]$ . Then,*

$$\mathbb{E}[f_\ell(\mathbf{x})x_{i^*}] \geq \mathbb{E}[f_\ell(\mathbf{x})x_j] \quad \text{for all } j \in [n].$$

**Lemma 2.4** (Useful properties of  $\mathcal{G}\text{-impurity}_f$ ).

1.  $\mathcal{G}\text{-impurity}_f(\text{empty tree}) = \mathcal{G}(\mathbb{E}[f]) \leq 1$ .
2.  $\text{dist}(f, T_f^\circ) \leq \mathcal{G}\text{-impurity}_f(T^\circ)$ .
3. For any leaf  $\ell$  of  $T^\circ$  and variable  $i \in [n]$ ,

$$\mathcal{G}\text{-purity-gain}_f(T^\circ, \ell, x_i) \geq 2^{-|\ell|} \cdot \frac{\kappa}{32} \cdot \text{Inf}_i(f_\ell)^2.$$

*Proof.* The first claim follows from the definition of  $\mathcal{G}\text{-impurity}_f$ . For the second claim, we have that

$$\begin{aligned} \text{dist}(f, T_f^\circ) &= \sum_{\text{leaves } \ell \in T^\circ} \left( \Pr[\mathbf{x} \text{ reaches } \ell] \cdot \text{bias}(f_\ell) \right) \\ &= \sum_{\text{leaves } \ell \in T^\circ} \left( 2^{-|\ell|} \cdot \text{bias}(f_\ell) \right) \\ &\leq \sum_{\text{leaves } \ell \in T^\circ} \left( 2^{-|\ell|} \cdot \mathcal{G}(\mathbb{E}[f_\ell]) \right) \quad (\text{Definition 2}) \\ &= \mathcal{G}\text{-impurity}_f(T^\circ). \end{aligned}$$

As for the third claim, we have that

$$\begin{aligned} & \mathcal{G}\text{-purity-gain}_f(T^\circ, \ell, x_i) \\ &= 2^{-|\ell|} \left( \mathcal{G}(\mathbb{E}[f_\ell]) - \mathbb{E}_{\mathbf{b} \sim \{\pm 1\}} [\mathcal{G}(\mathbb{E}[(f_\ell)_{x_i=b})]] \right) \left( \right) \\ &\geq 2^{-|\ell|} \cdot \frac{\kappa}{32} \cdot \text{Inf}_i(f_\ell)^2, \end{aligned}$$

where the final inequality is by Proposition 2.2.  $\square$

The following simple fact states that the error of the  $f$ -completion of a partial tree  $T^\circ$  cannot increase with further splits:

**Fact 2.5** (Splits cannot increase error). *Let  $T^\circ$  be a partial tree, and  $\tilde{T}^\circ$  be the partial tree that results from splitting a leaf of  $T^\circ$ . Then  $\text{dist}(f, \tilde{T}_f^\circ) \leq \text{dist}(f, T_f^\circ)$ .*

We are now ready to prove Theorem 2. By Fact 2.5, it suffices to show that error  $\text{opt}_s + \varepsilon$  is achieved after at most  $s^{O(\log s)}/\varepsilon^2$  iterations/splits. We do so by lower bounding the score of the leaf that is split by BUILDTOPDOWNDT in each iteration before error  $\text{opt}_s + \varepsilon$  is achieved. Let  $T^\circ$  be the size- $(j+1)$  partial tree that is built by BUILDTOPDOWNDT after  $j$  iterations. Suppose  $\text{dist}(f, T_f^\circ) > \text{opt}_s + \varepsilon$ , and let  $\ell^*$ ,  $x_{i^*}$  be the leaf and variable of  $T^\circ$  that is split in the  $(j+1)$ -st iteration. We claim that

$$\mathcal{G}\text{-purity-gain}_f(T^\circ, \ell^*, x_{i^*}) > \frac{\kappa \cdot \varepsilon^2}{32 \cdot (j+1)(\log s)^2}. \quad (1)$$

Writing  $x_{i(\ell)}$  to denote variable with the highest correlation with  $f_\ell$ , we have that

$$\begin{aligned} & \sum_{\text{leaves } \ell \in T^\circ} \left( \mathcal{G}\text{-purity-gain}_f(T^\circ, \ell, x_{i(\ell)}) \right) \\ &\geq \sum_{\text{leaves } \ell \in T^\circ} \left( 2^{-|\ell|} \cdot \frac{\kappa}{32} \cdot \text{Inf}_{i(\ell)}(f_\ell)^2 \right) \quad (\text{Lemma 2.4}) \\ &\geq \sum_{\text{leaves } \ell \in T^\circ} \left( 2^{-|\ell|} \cdot \frac{\kappa}{32} \cdot \left( \frac{\text{bias}(f_\ell) - \text{dist}(f_\ell, g_\ell)}{\text{size}(g_\ell)} \right)^2 \right) \quad (\text{Theorem 4}) \\ &\geq \frac{\kappa}{32} \cdot \sum_{\text{leaves } \ell \in T^\circ} \left( 2^{-|\ell|} \cdot \left( \frac{\text{bias}(f_\ell) - \text{dist}(f_\ell, g_\ell)}{\log s} \right)^2 \right) \quad (\text{size}(g_\ell) \leq \text{size}(g) = s) \\ &\geq \frac{\kappa}{32} \cdot \left[ \sum_{\text{leaves } \ell \in T^\circ} 2^{-|\ell|} \cdot \left( \frac{\text{bias}(f_\ell) - \text{dist}(f_\ell, g_\ell)}{\log s} \right) \right]^2 \quad (\text{Jensen's inequality}) \\ &= \frac{\kappa}{32} \cdot \frac{1}{(\log s)^2} \cdot \sum_{\text{leaves } \ell \in T^\circ} 2^{-|\ell|} \cdot \text{bias}(f_\ell) \\ &\quad - \sum_{\text{leaves } \ell \in T^\circ} \left( 2^{-|\ell|} \cdot \text{dist}(f_\ell, g_\ell) \right)^2 \\ &= \frac{\kappa}{32} \cdot \frac{1}{(\log s)^2} \cdot (\text{dist}(f, T_f^\circ) - \text{dist}(f, g))^2 \\ &> \frac{\kappa}{32} \cdot \frac{1}{(\log s)^2} \cdot ((\text{opt}_s + \varepsilon) - \text{opt}_s)^2 = \frac{\kappa}{32} \cdot \left( \frac{\varepsilon}{\log s} \right)^2. \end{aligned}$$

It follows that there must be at least one leaf and variable with purity gain greater than  $\kappa\varepsilon^2/32(j+1)(\log s)^2$ . Since BUILDTOPDOWNDT splits the leaf and variable with the largest purity gain, this establishes Equation (1).

Writing  $\tilde{T}^\circ$  to denote the partial tree that is obtained after BUILDTOPDOWNDT makes the single split with largest purity gain, we have that

$$\begin{aligned} & \mathcal{G}\text{-impurity}_f(\tilde{T}^\circ) \\ &= \mathcal{G}\text{-impurity}_f(T^\circ) - \mathcal{G}\text{-purity-gain}_f(T^\circ, \ell^*, x_{i^*}) \\ &\leq \mathcal{G}\text{-impurity}_f(T^\circ) - \frac{\kappa \cdot \varepsilon^2}{32 \cdot (j+1)(\log s)^2}. \end{aligned}$$

Combining this with the first and second claims of Lemma 2.4, we have the following: the value of the potential function starts off at at most 1 with  $T^\circ$  being the empty tree, decreases by at least  $\kappa \cdot \varepsilon^2 / 32j(\log s)^2$  with the  $j$ -th split, and error  $\text{opt}_s + \varepsilon$  is achieved once this value drops below  $\text{opt}_s + \varepsilon$ . Therefore, we can bound the number of splits necessary to ensure error  $\text{opt}_s + \varepsilon$  by the smallest  $t$  that satisfies:

$$\sum_{j=1}^t \left( \frac{\kappa \varepsilon^2}{32j(\log s)^2} \geq 1 - (\text{opt}_s + \varepsilon). \right)$$

Since

$$\sum_{j=1}^t \left( \frac{\kappa \varepsilon^2}{32j(\log s)^2} \geq \frac{\kappa \varepsilon^2 \log t}{32(\log s)^2}, \right)$$

we conclude that  $t \leq s^{O(\log s)/\kappa \varepsilon^2}$  suffices. This completes the proof of Theorem 2.

## 2.2. Trees for Real-Valued Features

We will prove an extension of Theorem 2 that applies to functions of real-valued features and arbitrary product distributions over inputs.

In order to state our result, we need to slightly restrict our definition of  $\text{opt}$  so it only considers ‘‘balanced’’ trees.

**Definition 8** (Balanced tree). *A decision tree  $T$  of depth  $d$  and size  $s$  is balanced if  $d = O(\log s)$ .*

**Definition 9** (balanced- $\text{opt}_s$ ). *For a function  $f : \mathbb{R}^n \rightarrow \{0, 1\}$ , product distribution  $\mathcal{D}$  over  $\mathbb{R}^n$ , and an integer  $s \in \mathbb{N}$ , we write  $\text{balanced\_opt}_{f, \mathcal{D}, s} \in [0, \frac{1}{2}]$  to denote the error of the best balanced size- $s$  decision tree for  $f$ :*

$$\text{balanced\_opt}_{f, \mathcal{D}, s} := \min \left\{ \Pr_{\mathbf{x} \sim \mathcal{D}} [T(\mathbf{x}) \neq f(\mathbf{x})] : \begin{array}{l} T \text{ is a balanced} \\ \text{size-}s \text{ decision tree} \end{array} \right\},$$

When  $f$  and  $\mathcal{D}$  are clear from context, we simply write  $\text{balanced\_opt}_s$ .

We can now formally state our main theorem.

**Theorem 1.** *Let  $f : \mathbb{R}^n \rightarrow \{0, 1\}$  be a monotone function and  $\mathcal{D}$  be any product distribution over  $\mathbb{R}^n$ . For any  $\kappa$ -strongly concave impurity heuristic  $\mathcal{G}$  and  $s \in \mathbb{N}$ ,*

$$\begin{aligned} & \text{TOPDOWNERROR}_{\mathcal{G}, \mathcal{D}}(f, s^{\tilde{O}((\log s)/\varepsilon^2)}) \\ &\leq \text{balanced\_opt}_s + \varepsilon. \end{aligned}$$

The proof of Theorem 1 will follow the same overall structure as our proof of Theorem 2. One key new ingredient is a generalization of Theorem 4 to real-valued features; this extension could be of independent interest:

**Theorem 5** (Extension of (Jain & Zhang, 2011) to trees for real-valued features). *Let  $f : \mathbb{R}^n \rightarrow \{0, 1\}$  be a monotone function,  $\mathcal{D}$  be an arbitrary product distribution over  $\mathbb{R}^n$ , and  $T$  be a size- $s$  balanced decision tree. Then, there exist  $i^* \in [n]$  and  $\theta^* \in \mathbb{R}$  for which  $\Pr_{\mathbf{x} \sim \mathcal{D}} [x_{i^*} \geq \theta^*] = \frac{1}{2}$  and*

$$\mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [f(\mathbf{x}) \cdot \mathbb{1}[x_{i^*} \geq \theta^*]] \not\geq \Omega \left( \frac{\varepsilon}{\log(s) \log \log(s/\varepsilon)} \right)$$

where  $\varepsilon := \text{bias}(f) - \text{dist}(f, T)$ .

where  $\text{bias}(f)$  and  $\text{dist}(f, T)$  are also measured with respect to  $\mathcal{D}$ .

To prove Theorem 1, we apply Theorem 5 in the same way Theorem 4 is used to prove Theorem 2. The full proof of Theorem 1 is deferred to the appendix.

## 3. Conclusion

We have given strengthened provable guarantees on the performance of widely employed and empirically successful top-down decision tree learning heuristics such as ID3, C4.5, and CART. Compared to previous works, our guarantees: (1) hold in the more realistic and challenging agnostic setting; (2) apply to all top-down heuristics and their associated impurity functions; (3) extend to the setting of real-valued features and arbitrary product distributions over the domain. Our main result shows that for all monotone functions  $f : \mathbb{R}^n \rightarrow \{0, 1\}$  and  $s \in \mathbb{N}$ , these top-down heuristics build a tree of size  $s^{\tilde{O}((\log s)/\varepsilon^2)}$  that achieves error within  $\varepsilon$  of that of the optimal balanced size- $s$  decision tree for  $f$ . We complement this with a near-matching lower bound. While our work was primarily motivated by the goal of understanding top-down heuristics, our results yield new guarantees that are not known to be achievable by any other algorithm, even ones that are not based on top-down heuristics.

There are several concrete avenues for future work:

1. Beyond monotonicity. As mentioned in the introduction, any top-down heuristic will fare badly on the parity functions  $f$ , in the sense of building a tree that is much larger than the optimal tree for  $f$ . Though broad and natural, the class of monotone functions is not the only class that excludes the parity function. Another fundamental property to consider is *noise stability* (see §2.4 of (O’Donnell, 2014)) — what guarantees can be made about the performance of these top-down heuristics when run on noise-stable functions?

2. Beyond product distributions. In this work our results hold for arbitrary product distributions over the domain, extending previous work that focuses on the uniform distribution. Could we establish provable distribution-independent guarantees, or failing that, perhaps provable guarantees for distributions with limited dependencies between coordinates?
3. Polynomial-size approximating trees. Our lower bound (Theorem 3) shows a monotone function such that any top-down heuristics has to build a tree of size  $s^{\tilde{\Omega}(\log s)}$  in order to achieve error  $\leq \text{opt}_s + \varepsilon$ . Results of (Blanc et al., 2020) show a similar lower bound of  $s^{\tilde{\Omega}(\sqrt[4]{\log s})}$  in the realizable setting. Are there broad and natural subclasses of monotone functions that evade these lower bounds, and for which polynomial size upper bounds do exist?

## Acknowledgements

We thank Michael Kim and the ICML reviewers for their helpful feedback and suggestions. LYT is supported by NSF grant CCF-192179 and NSF CAREER award CCF-1942123.

## References

- Blanc, G., Lange, J., and Tan, L.-Y. Top-down induction of decision trees: rigorous guarantees and inherent limitations. In *Proceedings of the 11th Innovations in Theoretical Computer Science Conference (ITCS)*, 2020. 1, 2, 1, 1.1, 1.1, 3, 2.1, 2.3, 3, A.1, C, 2, C.1, C, 6, C
- Blum, A., Furst, M., Jackson, J., Kearns, M., Mansour, Y., and Rudich, S. Weakly learning DNF and characterizing statistical query learning using Fourier analysis. In *Proceedings of the 26th Annual ACM Symposium on Theory of Computing (STOC)*, pp. 253–262, 1994. 1
- Blum, A., Burch, C., and Langford, J. On learning monotone boolean functions. In *Proceedings of the 39th Annual Symposium on Foundations of Computer Science (FOCS)*, pp. 408–415, 1998. 1
- Breiman, L. Random forests. *Machine learning*, 45(1): 5–32, 2001. 1
- Breiman, L. *Classification and regression trees*. Routledge, 2017. 1
- Brutzkus, A., Daniely, A., and Malach, E. On the Optimality of Trees Generated by ID3. *ArXiv*, abs/1907.05444, 2019a. 1.2
- Brutzkus, A., Daniely, A., and Malach, E. ID3 Learns Juntas for Smoothed Product Distributions. *ArXiv*, abs/1906.08654, 2019b. 1.2
- Bshouty, N. Exact learning via the monotone theory. *Information and Computation*, 123(1):146–153, 1995. 1
- Bshouty, N. and Tamon, C. On the Fourier spectrum of monotone functions. *Journal of the ACM*, 43(4):747–770, 1996. 1
- Chen, T. and Guestrin, C. Xgboost: A scalable tree boosting system. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 785–794, 2016. 1
- Dachman-Soled, D., Lee, H. K., Malkin, T., Servedio, R. A., Wan, A., and Wee, H. Optimal cryptographic hardness of learning monotone functions. *Theory of Computing*, 5(13):257–282, 2009. doi: 10.4086/toc.2009.v005a013. URL <http://www.theoryofcomputing.org/articles/v005a013>. 1
- Dachman-Soled, D., Feldman, V., Tan, L.-Y., Wan, A., and Wimmer, K. Approximate resilience, monotonicity, and the complexity of agnostic learning. In *Proceedings of the 26th Annual Symposium on Discrete Algorithms (SODA)*, pp. 498–511, 2015. 1
- Dietterich, T., Kearns, M., and Mansour, Y. Applying the weak learning framework to understand and improve C4.5. In *Proceedings of the 13th International Conference on Machine Learning (ICML)*, pp. 96–104, 1996. 1, 1.2
- Fiat, A. and Pechyony, D. Decision trees: More theoretical justification for practical algorithms. In *Proceedings of the 15th International Conference on Algorithmic Learning Theory (ALT)*, pp. 156–170, 2004. 1, 1.2
- Gopalan, P., Kalai, A., and Klivans, A. Agnostically learning decision trees. In *Proceedings of the 40th ACM Symposium on Theory of Computing (STOC)*, pp. 527–536, 2008. 1.2
- Hancock, T. and Mansour, Y. Learning monotone  $k$ - $\mu$  DNF formulas on product distributions. In *Proceedings of the 4th Annual Conference on Computational Learning Theory (COLT)*, pp. 179–193, 1991. 1
- Jackson, J., Lee, H., Servedio, R., and Wan, A. Learning Random Monotone DNF. *Discrete Applied Mathematics*, 159(5):259–271, 2011. 1
- Jain, R. and Zhang, S. The influence lower bound via query elimination. *Theory of Computing*, 7(1):147–153, 2011. 2, 5
- Kahn, J., Kalai, G., and Linial, N. The influence of variables on boolean functions. In *Proceedings of the 29th Annual Symposium on Foundations of Computer Science (FOCS)*, pp. 68–80, 1988. 2



- Kalai, A., Klivans, A., Mansour, Y., and Servedio, R. A. Agnostically learning halfspaces. *SIAM Journal on Computing*, 37(6):1777–1805, 2008. 1.2
- Kearns, M. Boosting theory towards practice: recent developments in decision tree induction and the weak learning framework (invited talk). In *Proceedings of the 13th National Conference on Artificial intelligence (AAAI)*, pp. 1337–1339, 1996. 1.2
- Kearns, M. and Mansour, Y. On the boosting ability of top-down decision tree learning algorithms. *Journal of Computer and System Sciences*, 58(1):109–128, 1999. 1, 1.2
- Kearns, M. and Valiant, L. Cryptographic limitations on learning Boolean formulae and finite automata. *Journal of the ACM*, 41(1):67–95, 1994. 1
- Kearns, M., Li, M., and Valiant, L. Learning Boolean formulas. *Journal of the ACM*, 41(6):1298–1328, 1994. 1
- Lee, H. *On the learnability of monotone functions*. PhD thesis, Columbia University, 2009. 1
- O’Donnell, R. *Analysis of Boolean Functions*. Cambridge University Press, 2014. Available at <http://analysisofbooleanfunctions.net/>. 1, A
- O’Donnell, R. and Servedio, R. Learning monotone decision trees in polynomial time. *SIAM Journal on Computing*, 37(3):827–844, 2007. 1, 1.2, C
- O’Donnell, R. and Wimmer, K. KKL, Kruskal–Katona, and Monotone Nets. *SIAM Journal on Computing*, 42(6):2375–2399, 2013. 1
- O’Donnell, R., Saks, M., Schramm, O., and Servedio, R. Every decision tree has an influential variable. In *Proceedings of the 46th Annual IEEE Symposium on Foundations of Computer Science (FOCS)*, pp. 31–39, 2005. 2, 2
- Pechyony, D. Decision trees: More theoretical justification for practical algorithms. Master’s thesis, Tel Aviv University, 2004. 1, 1.2
- Quinlan, R. Induction of decision trees. *Machine learning*, 1(1):81–106, 1986. 1
- Quinlan, R. *C4.5: Programs for Machine Learning*. Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, 1993. ISBN 1558602402. 1
- Sakai, Y. and Maruoka, A. Learning monotone log-term DNF formulas under the uniform distribution. *Theory of Computing Systems*, 33:17–33, 2000. 1
- Sellie, L. Learning random monotone DNF under the uniform distribution. In *Proceedings of the 21st Annual Conference on Learning Theory (COLT)*, pp. 181–192, 2008. 1
- Servedio, R. On learning monotone DNF under product distributions. *Information and Computation*, 193(1):57–74, 2004. 1
- Verbeurgt, K. Learning sub-classes of monotone DNF on the uniform distribution. In *Proceedings of the 9th Conference on Algorithmic Learning Theory (ALT)*, pp. 385–399, 1998. 1