

---

# Estimating decision tree learnability with polylogarithmic sample complexity

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 We show that top-down decision tree learning heuristics (such as ID3, C4.5, and  
2 CART) are amenable to highly efficient *learnability estimation*: for monotone target  
3 functions, the error of the decision tree hypothesis constructed by these heuristics  
4 can be estimated with *polylogarithmically* many labeled examples, exponentially  
5 smaller than the number necessary to run these heuristics, and indeed, exponentially  
6 smaller than information-theoretic minimum required to learn a good decision tree.  
7 This adds to a small but growing list of fundamental learning algorithms that have  
8 been shown to be amenable to learnability estimation.

9 En route to this result, we design and analyze sample-efficient *minibatch* versions  
10 of top-down decision tree learning heuristics and show that they achieve the same  
11 provable guarantees as the full-batch versions. We further give “active local”  
12 versions of these heuristics: given a test point  $x^*$ , we show how the label  $T(x^*)$   
13 of the decision tree hypothesis  $T$  can be computed with polylogarithmically many  
14 labeled examples, exponentially smaller than the number necessary to learn  $T$ .

## 15 1 Introduction

16 We study the problem of *estimating learnability*, recently introduced by Kong and Valiant [KV18]  
17 and Blum and Hu [BH18]. Consider a learning algorithm  $\mathcal{A}$  and a dataset  $S$  of *unlabeled* examples.  
18 Can we estimate the performance of  $\mathcal{A}$  on  $S$ —that is, the error of the hypothesis that  $\mathcal{A}$  would return  
19 if we were to label the entire dataset  $S$  and train  $\mathcal{A}$  on it—by labeling only very few of the examples  
20 in  $S$ ? Are there learning tasks and algorithms for which an accurate estimate of learnability can be  
21 obtained with far fewer labeled examples than the information-theoretic minimum required to learn a  
22 good hypothesis?

23 **Motivating applications.** Across domains and applications, the labeling of datasets is often an  
24 expensive process, requiring either significant computational resources or a large number of person-  
25 hours. There are therefore numerous natural scenarios in which an efficient learnability estimation  
26 procedure could serve as a useful exploratory precursor to learning. For example, suppose the error  
27 estimate returned by this procedure is large. This tells us that if we were to label the entire dataset  $S$   
28 and run  $\mathcal{A}$  on it, the error of the hypothesis  $h$  that  $\mathcal{A}$  would return is large. With this information,  
29 we may decide that  $h$  would not have been of much utility anyway, thereby saving ourselves the  
30 resources and effort to label the entire dataset  $S$  (and to run  $\mathcal{A}$ ). Alternatively, we may decide to  
31 collect more data or to enlarge the feature space of  $S$ , in hopes of improving the performance of  $\mathcal{A}$ .  
32 The learnability estimation procedure could again serve as a guide in this process, telling us *how*  
33 *much* the performance of  $\mathcal{A}$  would improve with these decisions. Relatedly, such a procedure could  
34 be useful for hyperparameter tuning, where the learning algorithm  $\mathcal{A}$  takes as input a parameter  $\rho$ ,  
35 and its performance improves with  $\rho$ , but its time and sample complexity also increases with  $\rho$ . The

36 learnability estimation procedure enables us to efficiently determine the best choice of  $\rho$  for our  
 37 application at hand, and run  $\mathcal{A}$  just a single time with this value of  $\rho$ . As a final example, such a  
 38 procedure could also be useful for dataset selection: given unlabeled training sets  $S_1, \dots, S_m$ , and  
 39 access to labeled examples from a test distribution  $\mathcal{D}$ , we can efficiently determine the  $S_i$  for which  $\mathcal{A}$   
 40 would produce a hypothesis that achieves the smallest error with respect to  $\mathcal{D}$ .

41 **Prior works on estimating learnability.** While this notion is still relatively new, there are by now  
 42 a number of works studying it in a variety of settings, including robust linear regression [KV18],  
 43 learning unions of intervals and  $k$ -Nearest-Neighbor algorithms [BH18], contextual bandits [KVB20],  
 44 learning Lipschitz functions, and the Nadaraya–Watson estimator in kernel regression [BBG20]. A  
 45 striking conceptual message has emerged from this line of work: it is often possible to estimate  
 46 learnability with far fewer labeled examples than the number required to run the corresponding  
 47 algorithm, and indeed, far fewer than the information-theoretic minimum required to learn a good  
 48 hypothesis.

## 49 1.1 Top-down decision tree learning

50 We study the problem of estimating learnability in the context of *decision tree learning*. Specifically,  
 51 we focus on *top-down* decision tree learning heuristics such as ID3 [Qui86], C4.5 [Qui93], and  
 52 CART [Bre17]. These classic and simple heuristics continue to be widely employed in everyday  
 53 machine learning applications and enjoy significant empirical success. They are also the core  
 54 subroutine in modern, state-of-the-art ensemble methods such as random forests [Bre01] and gradient  
 55 boosted trees [CG16].

56 We briefly describe how these top-down heuristics work, deferring the formal description to the  
 57 main body of this paper. Each such heuristic  $\text{TOPDOWN}_{\mathcal{G}}$  is defined by *impurity function*  $\mathcal{G} : [0, 1] \rightarrow [0, 1]$   
 58 which determines its splitting criterion.<sup>1</sup>  $\text{TOPDOWN}_{\mathcal{G}}$  takes as input a labeled dataset  
 59  $S \subseteq \mathcal{X} \times \{0, 1\}$  and a size parameter  $t \in \mathbb{N}$ , and constructs a size- $t$  decision tree for  $S$  in a *greedy*,  
 60 *top-down* fashion. It begins by querying  $\mathbb{1}[x_i \geq \theta]$  at the root of the tree, where  $x_i$  and  $\theta$  are chosen  
 61 to maximize the *purity gain with respect to*  $\mathcal{G}$ :

$$\mathcal{G}(\mathbb{E}[\mathbf{y}]) - (\Pr[\mathbf{x}_i \geq \theta] \cdot \mathcal{G}(\mathbb{E}[\mathbf{y} \mid \mathbf{x}_i \geq \theta]) + \Pr[\mathbf{x}_i < \theta] \cdot \mathcal{G}(\mathbb{E}[\mathbf{y} \mid \mathbf{x}_i < \theta])),$$

62 where the expectations and probabilities are with respect to  $(\mathbf{x}, \mathbf{y}) \sim S$ . More generally,  $\text{TOPDOWN}_{\mathcal{G}}$   
 63 grows its current tree  $T$  by splitting a leaf  $\ell \in T^\circ$  with a query to  $\mathbb{1}[x_i \geq \theta]$ , where  $\ell$ ,  $x_i$ , and  $\theta$  are  
 64 chosen to maximize:

$$\text{PurityGain}_{\mathcal{G}, S}(\ell, i, \theta) := \Pr[\mathbf{x} \text{ reaches } \ell] \cdot \text{LocalGain}_{\mathcal{G}, S}(\ell, i, \theta),$$

65 where

$$\begin{aligned} \text{LocalGain}_{\mathcal{G}, S}(\ell, i, \theta) &:= \mathcal{G}(\mathbb{E}[\mathbf{y} \mid \mathbf{x} \text{ reaches } \ell]) \\ &\quad - (\Pr[\mathbf{x}_i \geq \theta] \cdot \mathcal{G}(\mathbb{E}[\mathbf{y} \mid \mathbf{x} \text{ reaches } \ell, \mathbf{x}_i \geq \theta]) \\ &\quad \quad + \Pr[\mathbf{x}_i < \theta] \cdot \mathcal{G}(\mathbb{E}[\mathbf{y} \mid \mathbf{x} \text{ reaches } \ell, \mathbf{x}_i < \theta])). \end{aligned}$$

66 **Provable guarantees for monotone target functions [BLT20a].** Motivated by the tremendous  
 67 popularity and empirical successes of these top-down heuristics, there has been significant interest  
 68 and efforts in establishing provable guarantees on their performance [Kea96, DKM96, KM99, FP04,  
 69 Lee09, BDM19b, BDM19a, BLT20b, BLT20a]. The starting point of our work is a recent result of  
 70 Blanc et al. [BLT20a], which provides a guarantee on their performance when run on *monotone* target  
 71 functions, with respect to the uniform distribution:

72 **Theorem** (Theorem 2 of [BLT20a]). *Let  $f : \{\pm 1\}^d \rightarrow \{0, 1\}$  be a monotone target function*  
 73 *and  $\mathcal{G}$  be any impurity function. For  $s \in \mathbb{N}$  and  $\varepsilon, \delta \in (0, \frac{1}{2})$ , let  $t = s^{\Theta(\log s)/\varepsilon^2}$  and  $\mathbf{S}$  be*

<sup>1</sup>Impurity functions  $\mathcal{G}$  are restricted to be concave, symmetric around  $\frac{1}{2}$ , and to satisfy  $\mathcal{G}(0) = \mathcal{G}(1) = 0$   
 and  $\mathcal{G}(\frac{1}{2}) = 1$ . For example, ID3 and C4.5 use the binary entropy function  $\mathcal{G}(p) = H(p)$ , and the associated  
 purity gain is commonly referred to as information gain; CART uses the Gini criterion  $\mathcal{G}(p) = 4p(1 - p)$ ;  
 Kearns and Mansour proposed and analyzed the function  $\mathcal{G}(p) = 2\sqrt{p(1 - p)}$  [KM99]. The work of Dietterich,  
 Kearns, and Mansour [DKM96] provides a detailed discussion and experimental comparison of various impurity  
 functions.

74 a set of  $n$  labeled training examples  $(\mathbf{x}, f(\mathbf{x}))$  where  $\mathbf{x} \sim \{\pm 1\}^d$  is uniform random, and  $n =$   
 75  $\tilde{O}(t) \cdot \text{poly}(\log d, \log(1/\delta))$ .

76 With probability at least  $1 - \delta$  over the randomness of  $\mathcal{S}$ , the size- $t$  decision tree hypothesis constructed  
 77 by  $\text{TOPDOWN}_{\mathcal{G}}(t, \mathcal{S})$  satisfies  $\text{error}_f(T) := \Pr_{\mathbf{x} \sim \{\pm 1\}^d}[T(\mathbf{x}) \neq f(\mathbf{x})] \leq \text{opt}_s + \varepsilon$ , where  $\text{opt}_s$   
 78 denotes the error of the best size- $s$  decision tree for  $f$ .

79 We refer the reader to the introduction of [BLT20a] for a discussion of why assumptions on the  
 80 target function are necessary in order to establish provable guarantees. Briefly, as had been noted by  
 81 Kearns [Kea96], there are examples of simple non-monotone target functions  $f : \{\pm 1\}^d \rightarrow \{0, 1\}$ ,  
 82 computable by decision trees of constant size, for which any impurity-based heuristic may build a  
 83 complete tree of size  $\Omega(2^d)$  before achieving any non-trivial accuracy.

84 **Our contributions.** We give strengthened provable guarantees on the performance of top-down  
 85 decision tree learning heuristics, focusing on sample complexity. Our three main contributions are as  
 86 follows:

87 1. *Minibatch top-down decision tree learning.* We introduce and analyze  $\text{MINIBATCHTOPDOWN}_{\mathcal{G}}$ ,  
 88 a *minibatch* version of  $\text{TOPDOWN}_{\mathcal{G}}$  where the purity gain associated with each split is estimated  
 89 with only polylogarithmically many samples within the dataset  $S$  rather than all of  $S$ . For  
 90 all impurity functions  $\mathcal{G}$ , we show that  $\text{MINIBATCHTOPDOWN}_{\mathcal{G}}$  achieves the same provable  
 91 guarantees that those that [BLT20a] had established for the full-batch version  $\text{TOPDOWN}_{\mathcal{G}}$ .

92 2. *Active local learning.* We then study  $\text{MINIBATCHTOPDOWN}_{\mathcal{G}}$  within the recently-introduced  
 93 *active local learning* framework of Backurs, Blum, and Gupta [BBG20], and show that it admits  
 94 an efficient active local learner. Given active access to an *unlabeled* dataset  $S$  and a test point  $x^*$ ,  
 95 we show how  $T(x^*)$  can be computed by labeling only polylogarithmically many of the examples  
 96 in  $S$ , where  $T$  is the decision tree hypothesis that  $\text{MINIBATCHTOPDOWN}_{\mathcal{G}}$  would construct if  
 97 we were to label *all* of  $S$  and train  $\text{MINIBATCHTOPDOWN}_{\mathcal{G}}$  on it.

98 3. *Estimating learnability.* Building on both our results above, we show that  
 99  $\text{MINIBATCHTOPDOWN}_{\mathcal{G}}$  is amendable to highly-efficient learnability estimation. Given  
 100 active access to an unlabeled dataset  $S$ , we show that the error of  $T$  with respect to any test  
 101 distribution can be approximated by labeling only polylogarithmically many of the examples  
 102 in  $S$ , where  $T$  is the decision tree hypothesis that  $\text{MINIBATCHTOPDOWN}_{\mathcal{G}}$  would construct if  
 103 we were to label all of  $S$  and train  $\text{MINIBATCHTOPDOWN}_{\mathcal{G}}$  on it.

## 104 1.2 Formal statements of our results

105 **Feature space and distributional assumptions.** We work in the setting of binary attributes and  
 106 binary classification, i.e. we focus on the task of learning a target function  $f : \{\pm 1\}^d \rightarrow \{0, 1\}$ . We  
 107 will assume the learning algorithm receives uniform random examples  $\mathbf{x} \sim \{\pm 1\}^d$ , either labeled  
 108 or unlabeled. The error of a decision tree hypothesis  $T : \{\pm 1\}^d \rightarrow \{0, 1\}$  with respect to  $f$  is  
 109 defined to be  $\text{error}_f(T) := \Pr[f(\mathbf{x}) \neq T(\mathbf{x})]$  where  $\mathbf{x} \sim \{\pm 1\}^d$  is uniform random. We write  
 110  $\text{opt}_s(f)$  to denote  $\min\{\text{error}_f(T) : T \text{ is a size-}s \text{ decision tree}\}$ ; when  $f$  is clear from context we  
 111 simply write  $\text{opt}_s$ . We will also be interested in the error of  $T$  with respect to general test sets  
 112  $(\Pr_{(\mathbf{x}, \mathbf{y}) \sim S_{\text{test}}}[T(\mathbf{x}) \neq \mathbf{y}])$  and general test distributions  $(\Pr_{(\mathbf{x}, \mathbf{y}) \sim \mathcal{D}_{\text{test}}}[T(\mathbf{x}) \neq \mathbf{y}])$ .

113 **Notation and terminology.** For any decision tree  $T$ , we say the size of  $T$  is the number of leaves  
 114 in  $T$ . We refer to a decision tree with unlabeled leaves as a *partial tree*, and write  $T^\circ$  to denote such  
 115 trees. For a leaf  $\ell$  of a partial tree  $T^\circ$ , we write  $|\ell|$  to denote its depth within  $T^\circ$ , the number of  
 116 attributes queried along the path that leads to  $\ell$ . We say that an input  $x \in \{\pm 1\}^d$  is *consistent with*  
 117 *a leaf*  $\ell$  if  $x$  reaches  $\ell$  within  $T^\circ$ , and we write  $\ell_{T^\circ}(x)$  to denote the (unique) leaf  $\ell$  of  $T^\circ$  that  $x$   
 118 is consistent with. A function  $f : \{\pm 1\}^d \rightarrow \{0, 1\}$  is said to be *monotone* if for every coordinate  
 119  $i \in [d]$ , it is either non-decreasing with respect to  $i$  (i.e.  $f(x) \leq f(y)$  for all  $x, y \in \{\pm 1\}^d$  such  
 120 that  $x_i \leq y_i$ ) or non-increasing with respect to  $i$  (i.e.  $f(x) \geq f(y)$  for all  $x, y \in \{\pm 1\}^d$  such that  
 121  $x_i \leq y_i$ ).

122 We use **boldface** to denote random variables (e.g.  $\mathbf{x} \sim \{\pm 1\}^d$ ), and unless otherwise stated, all  
 123 probabilities and expectations are with respect to the uniform distribution. For  $p \in [0, 1]$ , we write

124  $\text{round}(p)$  to denote  $\mathbb{1}[p \geq \frac{1}{2}]$ . We reserve  $S$  to denote a labeled dataset and  $S^\circ$  to denote an unlabeled  
 125 dataset.

126 We are now ready to describe new algorithms and state our main results.

127 **Definition 1** (Minibatch). *Let  $S$  be a labeled dataset. A minibatch from  $S$ , denoted  $\mathbf{B} \sim \text{Batch}_b(S)$ ,  
 128 is a set of  $b$  uniform random points  $(x, y)$  chosen without replacement from  $S$ . More generally, for  
 129 a leaf  $\ell$ , a minibatch consistent with  $\ell$  from  $S$ , denoted  $\mathbf{B} \sim \text{Batch}_b(S, \ell)$ , is a set of  $b$  uniformly  
 130 random pairs chosen without replacement from among  $(x, y) \in S$  such that  $x$  is consistent with  $\ell$ . (In  
 131 both cases, if there are fewer than  $b$  such points, we return all of them.) Minibatches from unlabeled  
 132 datasets  $S^\circ$  are defined analogously.*

133 **Definition 2** (Minibatch completion of partial trees). *Given a partial tree  $T^\circ$ , we write  $T_{\text{Batch}_b(S)}^\circ$   
 134 to denote the tree obtained by labeling each leaf  $\ell \in T^\circ$  with  $\text{round}(\mathbb{E}_{(\mathbf{x}, f(\mathbf{x})) \sim \mathbf{B}}[f(\mathbf{x})])$  where  
 135  $\mathbf{B} \sim \text{Batch}_b(S, \ell)$ .*

MINIBATCHTOPDOWN $_{\mathcal{G}}(t, b, S)$ :

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

Define  $D := \log t + \log \log t$ .

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

1. **Score:** For each leaf  $\ell \in T^\circ$  of depth at most  $D$ , draw  $\mathbf{B} \sim \text{Batch}_b(S, \ell)$ . For each coordinate  $i \in [d]$ , compute:

$$\text{PurityGain}_{\mathcal{G}, \mathbf{B}}(\ell, i) := 2^{-|\ell|} \cdot \text{LocalGain}_{\mathcal{G}, \mathbf{B}}(\ell, i), \text{ where}$$

$$\begin{aligned} \text{LocalGain}_{\mathcal{G}, \mathbf{B}}(\ell, i) := & \mathcal{G}(\mathbb{E}[f(\mathbf{x})]) \\ & - \left( \frac{1}{2} \mathcal{G}(\mathbb{E}[f(\mathbf{x}) \mid \mathbf{x}_i = -1]) + \frac{1}{2} \mathcal{G}(\mathbb{E}[f(\mathbf{x}) \mid \mathbf{x}_i = 1]) \right), \end{aligned}$$

where the expectations are with respect to  $(\mathbf{x}, f(\mathbf{x})) \sim \mathbf{B}$ .

2. **Split:** Let  $(\ell^*, i^*)$  be the tuple that maximizes  $\text{PurityGain}_{\mathcal{G}, \mathbf{B}}(\ell, i)$ . Grow  $T^\circ$  by splitting  $\ell^*$  with a query to  $x_{i^*}$ .

Output  $T_{\text{Batch}_b(S)}^\circ$ .

Figure 1: MINIBATCHTOPDOWN $_{\mathcal{G}}$  takes as input a size parameter  $t$ , a minibatch size  $b$ , and a labeled dataset  $S$ . It outputs a size- $t$  decision tree hypothesis for  $f$ .

136 MINIBATCHTOPDOWN $_{\mathcal{G}}$  is a minibatch version of TOPDOWN $_{\mathcal{G}}$ , which we described informally  
 137 in Section 1.1 and include its full pseudocode in Appendix A. MINIBATCHTOPDOWN $_{\mathcal{G}}$  is more  
 138 efficient than TOPDOWN $_{\mathcal{G}}$  in two respects: first, purity gains and completions are computed with  
 139 respect to a minibatch  $\mathbf{B}$  of size  $b$  instead of all the entire dataset  $S$ ; second, MINIBATCHTOPDOWN $_{\mathcal{G}}$   
 140 never splits a leaf of depth greater than  $D$ , and hence constructs a decision tree of small size *and*  
 141 small depth, rather than just small size. (Looking ahead, both optimizations will be crucial for the  
 142 design of our sample-efficient active local learning and learnability estimation procedures.)

143 Our first result shows that MINIBATCHTOPDOWN $_{\mathcal{G}}$  achieves the same performance guarantees as  
 144 those that [BLT20a] had established for the full-batch version TOPDOWN $_{\mathcal{G}}$ :

145 **Theorem 1** (Provable guarantees for MINIBATCHTOPDOWN; informal version). *Let  $f : \{\pm 1\}^d \rightarrow$   
 146  $\{0, 1\}$  be a monotone target function and fix an impurity function  $\mathcal{G}$ . For any  $s \in \mathbb{N}$ ,  $\varepsilon, \delta \in (0, \frac{1}{2})$ ,  
 147 let  $t = s^{\Theta(\log s)/\varepsilon^2}$ , and  $\mathbf{S}$  be a set of  $n$  labeled training examples  $(\mathbf{x}, f(\mathbf{x}))$  where  $\mathbf{x} \sim \{\pm 1\}^d$  is  
 148 uniform random, and*

$$n = \tilde{O}(t) \cdot \text{poly}(\log d, \log(1/\delta)).$$

149 *If the minibatch size is at least*

$$b = \text{polylog}(t) \cdot \text{poly}(\log d, \log(1/\delta)),$$

150 then with probability at least  $1 - \delta$  over the randomness of  $\mathcal{S}$  and the draws of minibatches from  
 151 within  $\mathcal{S}$ , the size- $t$  decision tree hypothesis constructed by  $\text{MINIBATCHTOPDOWN}_{\mathcal{G}}(t, b, \mathcal{S})$  satis-  
 152 fies error  $f(T) \leq \text{opt}_s + \varepsilon$ .

153 Theorem 1 shows that it suffices for the minibatch size  $b$  of  $\text{MINIBATCHTOPDOWN}_{\mathcal{G}}$  to depend  
 154 polylogarithmically on  $t$ ; in contrast, the full-batch version  $\text{TOPDOWN}_{\mathcal{G}}$  uses the entire set  $S$  to  
 155 compute purity gains and determine its splits, and  $|S| = n$  has a superlinear dependence on  $t$ .

156 Our next algorithm is an implementation of  $\text{MINIBATCHTOPDOWN}_{\mathcal{G}}$  within the *active local learning*  
 157 framework of Backurs, Blum, and Gupta [BBG20]:

LOCALLEARNER $_{\mathcal{G}}(t, b, S^{\circ}, x^{\star})$ :

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

Define  $D := \log t + \log \log t$ .

Initialize  $e := 1$  and let  $\mathbf{B}_{\text{strands}}^{\circ}$  be  $b$  uniform random points from  $\{\pm 1\}^d$ .

while ( $e < t$ ) {

1. **Score:** For each leaf  $\ell \in \{\ell_{T^{\circ}}(x) : x \in \mathbf{B}_{\text{strands}}^{\circ} \cup \{x^{\star}\}\}$  of depth at most  $D$ , draw  $\mathbf{B}^{\circ} \sim \text{Batch}_b(S^{\circ}, \ell)$ , query  $f$ 's values on these points. For each coordinate  $i \in [d]$ , compute:
 
$$\text{PurityGain}_{\mathcal{G}, \mathbf{B}^{\circ}}(\ell, i) := 2^{-|\ell|} \cdot \text{LocalGain}_{\mathcal{G}, \mathbf{B}^{\circ}}(\ell, i), \text{ where}$$

$$\text{LocalGain}_{\mathcal{G}, \mathbf{B}^{\circ}}(\ell, i) := \mathcal{G}(\mathbb{E}[f(\mathbf{x})]) - \left( \frac{1}{2} \mathcal{G}(\mathbb{E}[f(\mathbf{x}) \mid \mathbf{x}_i = -1]) + \frac{1}{2} \mathcal{G}(\mathbb{E}[f(\mathbf{x}) \mid \mathbf{x}_i = 1]) \right),$$
 where the expectations are with respect to  $\mathbf{x} \sim \mathbf{B}^{\circ}$ .
2. **Split:** Let  $(\ell^{\star}, i^{\star})$  be the tuple that maximizes  $\text{PurityGain}_{\mathcal{G}, \mathbf{B}^{\circ}}(\ell, i)$ . Grow  $T^{\circ}$  by splitting  $\ell^{\star}$  with a query to  $x_{i^{\star}}$ .
3. **Estimate size:** Update our size estimate to
 
$$e = \mathbb{E}_{\mathbf{x} \sim \mathbf{B}_{\text{strands}}^{\circ}} [2^{|\ell_{T^{\circ}}(\mathbf{x})|}].$$

}

Draw  $\mathbf{B}^{\circ} \sim \text{Batch}_b(S^{\circ}, \ell_{T^{\circ}}(x^{\star}))$  and query  $f$ 's values on these points.

Output  $\text{round}(\mathbb{E}_{\mathbf{x} \sim \mathbf{B}^{\circ}}[f(\mathbf{x})])$ .

Figure 2: LOCALLEARNER $_{\mathcal{G}}$  takes as input a size parameter  $t$ , a minibatch size  $b$ , an unlabeled dataset  $S^{\circ}$ , and an input  $x^{\star}$ . It selectively queries  $f$ 's values on a few points within  $S^{\circ}$  and outputs  $T(x^{\star})$ , where  $T$  is a tree of size approximately  $t$  that  $\text{MINIBATCHTOPDOWN}_{\mathcal{G}}$  would return if we were to label all of  $S^{\circ}$  and train  $\text{MINIBATCHTOPDOWN}_{\mathcal{G}}$  on it.

158 **Theorem 2** (Active local version of  $\text{MINIBATCHTOPDOWN}$ ; informal version). Let  $f : \{\pm 1\}^d \rightarrow$   
 159  $\{0, 1\}$  be a target function,  $\mathcal{G}$  be an impurity function, and  $S^{\circ}$  be an unlabeled training set.

160 For all  $t \in \mathbb{N}$ ,  $\eta, \delta \in (0, \frac{1}{2})$ , if the minibatch size is at least  $b = \text{poly}(\log t, \log d, 1/\eta, \log(1/\delta))$ ,  
 161 then with probability at least  $1 - \delta$  over the randomness of  $\mathbf{B}_{\text{strands}}^{\circ}$ , we have that for all  $x^{\star} \in \{\pm 1\}^d$ ,  
 162 LOCALLEARNER $_{\mathcal{G}}(t, b, S^{\circ}, x^{\star})$  labels

$$q = O(b^2 \log t) = \text{polylog}(t) \cdot \text{poly}(\log d, 1/\eta, \log(1/\delta))$$

163 points within  $S^{\circ}$  and returns  $T(x^{\star})$ , where  $T$  is the size- $t'$  decision tree hypothesis that  
 164  $\text{MINIBATCHTOPDOWN}_{\mathcal{G}}(t', b, S)$  would construct,  $t' \in t(1 \pm \eta)$ , and  $S$  is the labeled dataset  
 165 obtained by labeling all of  $S^{\circ}$  with  $f$ 's values.<sup>2</sup>

<sup>2</sup>To ensure that LOCALLEARNER $_{\mathcal{G}}$  consistently labels all  $x^{\star}$  according to the same tree  $T$ , we run all invocations of LOCALLEARNER $_{\mathcal{G}}$  with the same outcomes of randomness for  $\mathbf{B}_{\text{strands}}^{\circ}$  and draws of minibatches. Similarly, if one then wished to actually construct this tree  $T$ , they would run  $\text{MINIBATCHTOPDOWN}_{\mathcal{G}}$  with these same outcomes of randomness.

166 Theorem 2 yields, as a fairly straightforward consequence, our learnability estimation procedure  
 167  $\text{Est}_{\mathcal{G}}$  that estimates the performance of  $\text{MINIBATCHTOPDOWN}_{\mathcal{G}}$  with respect to any test set  $S_{\text{test}}$ :

168 **Theorem 3** (Estimating learnability of  $\text{MINIBATCHTOPDOWN}$ ; informal version). *Let  $f : \{\pm 1\}^d \rightarrow$   
 169  $\{0, 1\}$  be a target function,  $\mathcal{G}$  be an impurity function,  $S^\circ$  be an unlabeled training set, and  $S_{\text{test}}$  be  
 170 a labeled test set.*

171 *For all  $t \in \mathbb{N}$  and  $\eta, \delta \in (0, \frac{1}{2})$ , if the minibatch size  $b$  is as in Theorem 2, then with probability at  
 172 least  $1 - \delta$  over the randomness of the draws of minibatches from within  $S^\circ$ ,  $\text{EST}_{\mathcal{G}}(t, b, S^\circ, S_{\text{test}})$   
 173 labels*

$$q = O(|S_{\text{test}}| \cdot b \log t + b^2 \log t) = |S_{\text{test}}| \cdot \text{polylog}(t) \cdot \text{poly}(\log d, 1/\eta, \log(1/\delta))$$

174 *points within  $S^\circ$  and returns the error of  $T$  with respect to  $S_{\text{test}}$ ,*

$$\text{error}_{S_{\text{test}}}(T) := \Pr_{(\mathbf{x}, \mathbf{y}) \sim S_{\text{test}}} [T(\mathbf{x}) \neq \mathbf{y}],$$

175 *where  $T$  is as in Theorem 2.*

176 We remark that the labels of the test set  $S_{\text{test}}$  in Theorem 3 need not be consistent with  $f$ . Indeed, as  
 177 an example application of Theorem 3, we can let  $S_{\text{test}}$  be  $\Theta(\log(1/\delta)/\varepsilon^2)$  many labeled examples  
 178  $(\mathbf{x}, \mathbf{y})$  drawn from an arbitrary test distribution  $\mathcal{D}_{\text{test}}$  over  $\{\pm 1\}^d \times \{0, 1\}$ , where the marginal over  
 179  $\{\pm 1\}^d$  need not be uniform and the labels need not be consistent with  $f$ , and with probability at  
 180 least  $1 - \delta$ , the output of  $\text{Est}_{\mathcal{G}}$  will be within  $\pm \varepsilon$  of  $\Pr_{(\mathbf{x}, \mathbf{y}) \sim \mathcal{D}_{\text{test}}} [T(\mathbf{x}) \neq \mathbf{y}]$ .

## 181 2 Proof overview for Theorem 1

182 Our proof of Theorem 1 builds upon and extends the analysis in [BLT20a]. (Recall that [BLT20a]  
 183 analyzed the full-batch version  $\text{TOPDOWN}_{\mathcal{G}}$ , which we have included in Appendix A of this paper,  
 184 and their guarantee concerning its performance is their Theorem 2, which we have stated in Section 1.1  
 185 of this paper). In this section we give a high-level overview of both [BLT20a]’s and our proof strategy,  
 186 in tandem with a description of the technical challenges that arise as we try to strengthen [BLT20a]’s  
 187 Theorem 2 to our Theorem 1.

188 Let  $f : \{\pm 1\}^d \rightarrow \{0, 1\}$  be a monotone function and fix an impurity function  $\mathcal{G}$ . Let  $T^\circ$  be  
 189 a partial tree that is being built by either  $\text{TOPDOWN}_{\mathcal{G}}$  or  $\text{MINIBATCHTOPDOWN}_{\mathcal{G}}$ . Recall that  
 190  $\text{TOPDOWN}_{\mathcal{G}}$  and  $\text{MINIBATCHTOPDOWN}_{\mathcal{G}}$  compute, for each leaf  $\ell \in T^\circ$  and coordinate  $i \in [d]$ ,  
 191  $\text{PurityGain}_{\mathcal{G}, S}(\ell, i)$  and  $\text{PurityGain}_{\mathcal{G}, B}(\ell, i)$  respectively. Both these quantities can be thought of  
 192 as estimates of the *true* purity gain:

$$\begin{aligned} \text{PurityGain}_{\mathcal{G}, f}(\ell, i) &:= 2^{-|\ell|} \cdot \text{LocalGain}_{\mathcal{G}, f}(\ell, i) \text{ where} \\ \text{LocalGain}_{\mathcal{G}, f}(\ell, i) &:= \mathcal{G}(\mathbb{E}[f(\mathbf{x}) \mid \mathbf{x} \text{ reaches } \ell]) \\ &\quad - \left(\frac{1}{2} \mathcal{G}(\mathbb{E}[f(\mathbf{x}) \mid \mathbf{x} \text{ reaches } \ell, x_i = -1])\right) \\ &\quad + \frac{1}{2} \mathcal{G}(\mathbb{E}[f(\mathbf{x}) \mid \mathbf{x} \text{ reaches } \ell, x_i = 1]), \end{aligned}$$

193 where here and throughout this section, all expectations are with respect to a uniform random  
 194  $\mathbf{x} \sim \{\pm 1\}^d$ . The fact that  $\text{MINIBATCHTOPDOWN}_{\mathcal{G}}$ ’s estimates of this true purity gain are based on  
 195 minibatches  $B$  of size exponentially smaller than that of the full sample set  $S$ —and hence could be ex-  
 196 pectationally less accurate—is a major source of technical challenges that arise in extending [BLT20a]’s  
 197 guarantees for  $\text{TOPDOWN}_{\mathcal{G}}$  to  $\text{MINIBATCHTOPDOWN}_{\mathcal{G}}$ .

198 [BLT20a] considers the potential function:

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

199 The following fact about this potential function  $\mathcal{G}\text{-impurity}_f$  is straightforward to verify (and is  
 200 proved in [BLT20a]):

201 **Fact 2.1.** *For any partial tree  $T^\circ$ , leaf  $\ell \in T^\circ$ , and coordinate  $i \in [d]$ , let  $\tilde{T}^\circ$  be the tree obtained  
 202 from  $T^\circ$  by splitting  $\ell$  with a query to  $x_i$ . Then,*

$$\mathcal{G}\text{-impurity}_f(\tilde{T}^\circ) = \mathcal{G}\text{-impurity}_f(T^\circ) - \text{PurityGain}_{\mathcal{G}, f}(\ell, i).$$

203 A key ingredient in [BLT20a]’s analysis is a proof that as long as  $\text{error}_f(T_S^\circ) > \text{opt}_s + \varepsilon$  (where  
 204  $T_S^\circ$  denotes the completion of  $T^\circ$  with respect to the full batch  $S$ ; see Appendix A), there must  
 205 be a leaf  $\ell \in T^\circ$  and coordinate  $i$  with high true purity gain,  $\text{PurityGain}_{\mathcal{G},f}(\ell, i) \geq \text{poly}(\varepsilon/t)$ .  
 206 Since  $\text{TOPDOWN}_{\mathcal{G}}$ ’s estimates  $\text{PurityGain}_{\mathcal{G},S}$  of  $\text{PurityGain}_{\mathcal{G},f}$  are with respect to a sample of  
 207 size  $|S| \geq \text{poly}(t/\varepsilon)$ , it follows that  $\text{TOPDOWN}_{\mathcal{G}}$  will make a split for which the true purity gain is  
 208 indeed  $\text{poly}(\varepsilon/t)$ . By Fact 2.1, such a split constitutes good progress with respect to the potential  
 209 function  $\mathcal{G}$ -impurity $_f$ . Summarizing, [BLT20a] that shows until  $\text{error}_f(T_S^\circ) < \text{opt}_s + \varepsilon$  is achieved,  
 210 every split that  $\text{TOPDOWN}_{\mathcal{G}}$  makes has high true purity gain, and hence constitutes good progress  
 211 with respect to the potential function  $\mathcal{G}$ -impurity $_f$ .

212 The key technical difficulty in analyzing  $\text{MINIBATCHTOPDOWN}_{\mathcal{G}}$  instead of  $\text{TOPDOWN}_{\mathcal{G}}$  is that  
 213  $\text{MINIBATCHTOPDOWN}_{\mathcal{G}}$  is not guaranteed to choose a split with high true purity gain: it could make  
 214 splits for which its estimate  $\text{PurityGain}_{\mathcal{G},B}(\ell, i)$  is high, but the true purity gain  $\text{PurityGain}_{\mathcal{G},f}(\ell, i)$   
 215 is actually tiny. In fact, unless we use batches of size  $b \geq \text{poly}(t)$ , exponentially larger than the  
 216  $b = \text{polylog}(t)$  of Theorem 1,  $\text{MINIBATCHTOPDOWN}_{\mathcal{G}}$  could make splits that result in zero true  
 217 purity gain, and hence constitute zero progress with respect to the potential function  $\mathcal{G}$ -impurity $_f$ .

218 To overcome this challenge, we instead show that *most* splits  $\text{MINIBATCHTOPDOWN}_{\mathcal{G}}$  makes have  
 219 high true purity gain. We first show that with high probability over the draws of minibatches  $B$ , if  
 220  $\text{MINIBATCHTOPDOWN}_{\mathcal{G}}$  splits a leaf that is neither too shallow nor too deep within  $T^\circ$ , then this  
 221 split has high true purity gain (Lemma B.5). We then show the following two lemmas:

- 222 1. Lemma B.6: If  $\text{MINIBATCHTOPDOWN}_{\mathcal{G}}$  splits a leaf of  $T^\circ$  that is sufficiently deep, then it must  
 223 be the case that  $\text{error}_f(T_{\text{Batch}_b(S)}^\circ) \leq \text{opt}_s + \varepsilon$ , i.e. the current tree already achieves sufficiently  
 224 small error. With this Lemma, we are able to define  $\text{MINIBATCHTOPDOWN}_{\mathcal{G}}$  to never split a  
 225 leaf that is too deep, while retaining guarantees on its performance.
- 226 2. Lemma B.7: This lemma shows that only a small fraction of splits made by  
 227  $\text{MINIBATCHTOPDOWN}_{\mathcal{G}}$  can be too shallow.

228 Combining the above Lemmas, we are able to prove Theorem 1. We defer the proof to Appendix B.

### 229 3 Proof overviews for Theorems 2 and 3

230 We begin with a proof overview for Theorem 2. Let  $T$  be the decision tree hypothesis that  
 231  $\text{MINIBATCHTOPDOWN}_{\mathcal{G}}$  would construct if we were to all of  $S^\circ$  and train  $\text{MINIBATCHTOPDOWN}_{\mathcal{G}}$   
 232 on it. Our goal is to efficiently compute  $T(x^*)$  for a given  $x^*$  by selectively labeling only  $q$  points  
 233 within  $S^\circ$ , where  $q$  is exponentially smaller than the sample complexity of learning and construct-  
 234 ing  $T$ .

235 Intuitively, we would like  $\text{LOCALLEARNER}_{\mathcal{G}}$  to only grow the single “strand” within  $T$  required to  
 236 compute  $T(x^*)$  instead of the entire tree  $T$ —this “strand” is simply the root-to-leaf path of  $T$  that  $x^*$   
 237 follows. The key challenge that arises in implementing this plan is: how does  $\text{LOCALLEARNER}_{\mathcal{G}}$   
 238 know when to terminate this strand (i.e. how does it know when it has reached a leaf of  $T$ )?  
 239  $\text{MINIBATCHTOPDOWN}_{\mathcal{G}}$ , the “global” algorithm that  $\text{LOCALLEARNER}_{\mathcal{G}}$  is trying to simulate,  
 240 terminates when the tree is of size  $t$ . As  $\text{LOCALLEARNER}_{\mathcal{G}}$  grows the strand corresponding to  $x^*$ ,  
 241 how could it estimate the size of the overall tree without actually growing it? In other words, it is not  
 242 clear how one would define the stopping criterion of the while loop in the following pseudocode:

```

Initialize  $\ell$  to be the leaf of the empty tree.
while (stopping criterion) {
  1. Draw  $B^\circ \sim \text{Batch}_b(S^\circ, \ell)$  and query  $f$ ’s values on these points. Let  $i^*$  be the
     coordinate that maximizes  $\text{PurityGain}_{\mathcal{G},B^\circ}(\ell, i)$  among all  $i \in [d]$ .
  2. Extend  $\ell$  according to the value of  $x_{i^*}^*$ .
}
Draw  $B^\circ \sim \text{Batch}_b(S, \ell)$  and query  $f$ ’s values on these points.
Output  $\text{round}(\mathbb{E}_{x \sim B^\circ}[f(x)])$ .

```

243 Roughly speaking, we want “*stopping criterion*” to answer the following question: if we grew a  
 244 size- $t$  tree using  $\text{MINIBATCHTOPDOWN}_{\mathcal{G}}$  (on the labeled version of  $S^\circ$ ), would  $\ell$  be a leaf of the  
 245 resulting tree, or would it be an internal node? Nearly equivalently, with access to just a single strand  
 246 of a tree, we wish to estimate the size of that tree. If that size is  $t$ , then we stop the while loop.

247 It is not possible to accurately estimate the size of a tree using just a single strand. However, by comput-  
 248 ing a small number of random strands, we can get an accurate size estimator. In Appendix C, we show  
 249 that for  $\mathbf{x}_1, \dots, \mathbf{x}_m$  chosen uniformly at random from  $\{\pm 1\}^d$ , the estimator  $e := \frac{1}{m} \sum_{i=1}^m 2^{|\ell_T(\mathbf{x}_i)|}$   
 250 accurately estimates the size of  $T$ , as long as the depth of  $T$  is not too large. Therefore, rather than  
 251 growing only the root-to-leaf path for  $x^*$ ,  $\text{LOCALLEARNER}_{\mathcal{G}}$  samples random additional inputs,  
 252  $\mathbf{x}_1, \dots, \mathbf{x}_m$ . Then, it simultaneously grows the strands for the root-to-leaf paths of  $x^*$  as well as  
 253  $\mathbf{x}_1, \dots, \mathbf{x}_m$ . These strands do not all grow at the same “rate”, as we want  $\text{LOCALLEARNER}_{\mathcal{G}}$   
 254 to make splits in the same order as  $\text{MINIBATCHTOPDOWN}_{\mathcal{G}}$  does. As long as it does this, we can use  
 255 the size estimator to, at any step, accurately estimate the size of tree  $\text{MINIBATCHTOPDOWN}_{\mathcal{G}}$  would  
 256 need to build for all the current strands to end at leaves.  $\text{LOCALLEARNER}_{\mathcal{G}}$  terminates when its  
 257 estimate of this size is  $t$ .

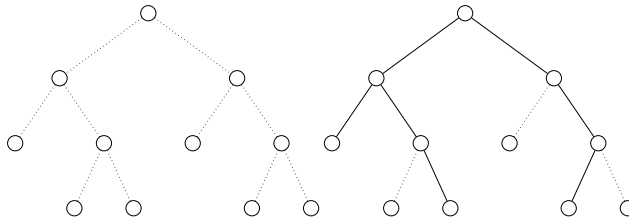


Figure 3: Rather than growing the entire tree  $T$  (depicted on the LHS) as  $\text{MINIBATCHTOPDOWN}_{\mathcal{G}}$  does,  $\text{LOCALLEARNER}_{\mathcal{G}}$  only grows  $m + 1$  strands within  $T$  (depicted on the RHS), corresponding to the given input  $x^*$  and  $m$  additional random inputs  $\mathbf{x}_1, \dots, \mathbf{x}_m \sim \{\pm 1\}^d$ .

258 We back the above intuition for  $\text{LOCALLEARNER}_{\mathcal{G}}$  with proofs. In Appendix D, we show that the  
 259 output of  $\text{LOCALLEARNER}_{\mathcal{G}}$  for size parameter  $t$  is  $T(x^*)$ , where  $T$  is size- $t'$  tree produced by  
 260  $\text{MINIBATCHTOPDOWN}_{\mathcal{G}}$  where  $t \in t'(1 \pm \eta)$ . We also show that  $\text{LOCALLEARNER}_{\mathcal{G}}$  needs to only  
 261 label polylogarithmic many points within  $S^\circ$  to compute  $T(x^*)$ . This completes our proof overview  
 262 for Theorem 2, and Theorem 3 is a straightforward consequence of Theorem 2.

## 263 4 Conclusion

264 We have given strengthened provable guarantees on the performance of popular and empirically  
 265 successful top-down decision tree learning heuristics such as ID3, C4.5, and CART, focusing  
 266 on sample complexity. First, we designed and analyzed *minibatch* versions of these heuristics,  
 267  $\text{MINIBATCHTOPDOWN}_{\mathcal{G}}$ , and proved that they achieve the same performance guarantees as the  
 268 full-batch versions. We then gave an implementation of  $\text{MINIBATCHTOPDOWN}_{\mathcal{G}}$  within the recently-  
 269 introduced active local learning framework of [BBG20]. Building on these results, we showed that  
 270  $\text{MINIBATCHTOPDOWN}_{\mathcal{G}}$  is amenable to highly efficient learnability estimation [KV18, BH18]: its  
 271 performance can be estimated accurately by selectively labeling very few examples.

272 As discussed in [KV18, BH18], this new notion of learnability estimation opens up a whole host of  
 273 theoretical and empirical directions for future work. We discuss several concrete ones that are most  
 274 relevant to our work. Our algorithm  $\text{Est}_{\mathcal{G}}$  efficiently and accurately estimates the quality, relative to  
 275 a test set  $S_{\text{test}}$ , of the hypothesis that  $\text{MINIBATCHTOPDOWN}_{\mathcal{G}}$  would produce if trained on a set  $S^\circ$ .  
 276 Could  $\text{Est}_{\mathcal{G}}$  be more broadly useful in assessing the quality of the *training data*  $S^\circ$  itself, relative to  
 277  $S_{\text{test}}$ ? Could its estimates provide guarantees on the performance of other algorithms when trained  
 278 in  $S^\circ$  and tested on  $S_{\text{test}}$ ? It would also be interesting to explore applications of our algorithms to  
 279 the design of training sets. Given training sets  $S_1, \dots, S_m$ ,  $\text{Est}_{\mathcal{G}}$  allows us to efficiently determine  
 280 the  $S_i$  for which  $\text{MINIBATCHTOPDOWN}_{\mathcal{G}}$  would produce a hypothesis that achieves the smallest  
 281 error with respect to  $S_{\text{test}}$ . Could  $\text{Est}_{\mathcal{G}}$  or extensions of it be useful in efficiently *creating an*  $S^*$ ,  
 282 comprising data from each  $S_i$ , that is of higher quality than any  $S_i$  individually? Finally, while we  
 283 have focused on top-down heuristics for learning a single decision tree in this work, a natural next  
 284 step would be to design and analyze learnability estimation procedures for ensemble methods such as  
 285 random forests and gradient boosted trees.



286 **Broader Impact**

287 This work does not present any foreseeable societal consequence.

288 **References**

- 289 [BBG20] Arturs Backurs, Avrim Blum, and Neha Gupta. Active local learning. In *Proceedings of*  
290 *the 33rd Conference On Learning Theory (COLT)*. Proceedings of Machine Learning  
291 Research, 2020.
- 292 [BDM19a] Alon Brutzkus, Amit Daniely, and Eran Malach. ID3 Learns Juntas for Smoothed  
293 Product Distributions. *ArXiv*, abs/1906.08654, 2019.
- 294 [BDM19b] Alon Brutzkus, Amit Daniely, and Eran Malach. On the Optimality of Trees Generated  
295 by ID3. *ArXiv*, abs/1907.05444, 2019.
- 296 [BH18] Avrim Blum and Lunjia Hu. Active tolerant testing. In *Proceedings of the 31st Conference*  
297 *On Learning Theory (COLT)*, volume 75, pages 474–497. Proceedings of Machine  
298 Learning Research, 2018.
- 299 [BLT20a] Guy Blanc, Jane Lange, and Li-Yang Tan. Provable guarantees for decision tree induction:  
300 the agnostic setting. In *Proceedings of the 37th International Conference on Machine*  
301 *Learning (ICML)*, 2020. Available at <https://arxiv.org/abs/2006.00743>.
- 302 [BLT20b] Guy Blanc, Jane Lange, and Li-Yang Tan. Top-down induction of decision trees: rigorous  
303 guarantees and inherent limitations. In *Proceedings of the 11th Innovations in Theoretical*  
304 *Computer Science Conference (ITCS)*, volume 151, pages 1–44, 2020.
- 305 [Bre01] Leo Breiman. Random forests. *Machine learning*, 45(1):5–32, 2001.
- 306 [Bre17] Leo Breiman. *Classification and regression trees*. Routledge, 2017.
- 307 [CG16] Tianqi Chen and Carlos Guestrin. Xgboost: A scalable tree boosting system. In *Proceed-*  
308 *ings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and*  
309 *Data Mining*, pages 785–794, 2016.
- 310 [DKM96] Tom Dietterich, Michael Kearns, and Yishay Mansour. Applying the weak learning  
311 framework to understand and improve C4.5. In *Proceedings of the 13th International*  
312 *Conference on Machine Learning (ICML)*, pages 96–104, 1996.
- 313 [FP04] Amos Fiat and Dmitry Pechyony. Decision trees: More theoretical justification for  
314 practical algorithms. In *Proceedings of the 15th International Conference on Algorithmic*  
315 *Learning Theory (ALT)*, pages 156–170, 2004.
- 316 [Kea96] Michael Kearns. Boosting theory towards practice: recent developments in decision tree  
317 induction and the weak learning framework (invited talk). In *Proceedings of the 13th*  
318 *National Conference on Artificial intelligence (AAAI)*, pages 1337–1339, 1996.
- 319 [KM99] Michael Kearns and Yishay Mansour. On the boosting ability of top-down decision tree  
320 learning algorithms. *Journal of Computer and System Sciences*, 58(1):109–128, 1999.
- 321 [KV18] Weihao Kong and Gregory Valiant. Estimating learnability in the sublinear data regime.  
322 In *31st Annual Conference on Neural Information Processing Systems (NeurIPS)*, pages  
323 5460–5469, 2018.
- 324 [KVB20] Weihao Kong, Gregory Valiant, and Emma Brunskill. Sublinear optimal policy value  
325 estimation in contextual bandits. In *Proceedings of the 23rd International Conference*  
326 *on Artificial Intelligence and Statistics (AISTATS)*, 2020.
- 327 [Lee09] Homin Lee. *On the learnability of monotone functions*. PhD thesis, Columbia University,  
328 2009.
- 329 [Qui86] Ross Quinlan. Induction of decision trees. *Machine learning*, 1(1):81–106, 1986.
- 330 [Qui93] Ross Quinlan. *C4.5: Programs for Machine Learning*. Morgan Kaufmann Publishers  
331 Inc., San Francisco, CA, USA, 1993.

332 **A Pseudocode for TOPDOWN $\mathcal{G}$**

TOPDOWN $\mathcal{G}(t, S)$ :  
 Initialize  $T^\circ$  to be the empty tree.  
 while ( $\text{size}(T^\circ) < t$ ) {  
   1. **Score:** For each leaf  $\ell \in T^\circ$  and coordinate  $i \in [d]$ , compute:  
     PurityGain $\mathcal{G}, S(\ell, i) := 2^{-|\ell|} \cdot \text{LocalGain}_{\mathcal{G}, S}(\ell, i)$ , where  
     LocalGain $\mathcal{G}, S(\ell, i) := \mathcal{G}(\mathbb{E}[f(\mathbf{x}) \mid \mathbf{x} \text{ reaches } \ell])$   
      $- (\frac{1}{2} \cdot \mathcal{G}(\mathbb{E}[f(\mathbf{x}) \mid \mathbf{x} \text{ reaches } \ell, \mathbf{x}_i = -1])$   
      $+ \frac{1}{2} \cdot \mathcal{G}(\mathbb{E}[f(\mathbf{x}) \mid \mathbf{x} \text{ reaches } \ell, \mathbf{x}_i = 1]))$ ,  
     where  $(\mathbf{x}, f(\mathbf{x})) \sim S$ .  
   2. **Split:** Let  $(\ell^*, i^*)$  be the tuple that maximizes PurityGain $\mathcal{G}, S(\ell, i)$ . Grow  $T^\circ$  by  
     splitting  $\ell^*$  with a query to  $x_{i^*}$ .  
 }  
 Output  $T_S^\circ$ , the completion of  $T^\circ$  with respect to  $S$ : label each leaf  $\ell \in T^\circ$  with  
 round( $\mathbb{E}[f(\mathbf{x}) \mid \mathbf{x} \text{ reaches } \ell]$ ), where  $(\mathbf{x}, f(\mathbf{x})) \sim S$ .

333 **B Provable guarantees for MINIBATCHTOPDOWN**

334 In this section we will prove Theorem 1. We first need a couple of definitions:

335 **Definition 3** (Hölder continuous). For  $C, \alpha > 0$ , an impurity function  $\mathcal{G} : [0, 1] \rightarrow [0, 1]$  is  $(C, \alpha)$ -  
 336 Hölder continuous if, for all  $a, b \in [0, 1]$ ,

$$|\mathcal{G}(a) - \mathcal{G}(b)| \leq C|a - b|^\alpha.$$

337 **Definition 4** (Strong concavity). For  $\kappa > 0$ , an impurity function  $\mathcal{G} : [0, 1] \rightarrow [0, 1]$  is  $\kappa$ -strongly  
 338 concave if for all  $a, b \in [0, 1]$ ,

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

339 **Theorem 4** (Provable guarantee for MINIBATCHTOPDOWN; formal version of Theorem 1). Let  
 340  $f : \{\pm 1\}^d \rightarrow \{0, 1\}$  be a monotone target function and  $\mathcal{G}$  be any  $\kappa$ -strongly concave and  $(C, \alpha)$ -  
 341 Hölder continuous impurity function. For any  $s \in \mathbb{N}$ ,  $\varepsilon, \delta \in (0, \frac{1}{2})$ , let  $t = s^{\Theta(\log(s))/\varepsilon^2}$ , and  $\mathbf{S}$  be a  
 342 set of  $n$  labeled training examples  $(\mathbf{x}, f(\mathbf{x}))$  where  $\mathbf{x} \sim \{\pm 1\}^d$  is uniform random, and

$$n = t \cdot \Omega\left(\left(\frac{C^2 \log(s)^4}{\kappa^2 \varepsilon^4}\right)^{\frac{1}{\alpha}} \cdot \log\left(\frac{td}{\delta}\right) \cdot \log t\right).$$

343 If the minibatch size is at least

$$b = \Omega\left(\left(\frac{C^2 \log(s)^4}{\kappa^2 \varepsilon^4}\right)^{\frac{1}{\alpha}} \cdot \log\left(\frac{td}{\delta}\right)\right),$$

344 then with probability at least  $1 - \delta$  over the randomness of  $\mathbf{S}$  and the draws of minibatches from within  
 345  $\mathbf{S}$ , the size- $t$  decision tree hypothesis constructed by MINIBATCHTOPDOWN $\mathcal{G}(t, b, \mathbf{S})$  satisfies  
 346 error $_f(T) \leq \text{opt}_s + \varepsilon$ .

347 **B.1 Properties of batches**

348 We begin by specifying how large the batch size has to be for accurate estimates of local gain. Later  
 349 on, we will turn accurate estimates of local gain to estimates of purity gain that are accurate at least  
 350 half the time.

351 **Lemma B.1** (Every leaf has a batch of size  $b_{\min}$ ). *Let*

$$b_{\min} = \max \left( 8, 2 \cdot \left( \frac{2C}{\Delta} \right)^{\frac{2}{\alpha}} \right) \cdot \log_e \left( \frac{9td}{\delta} \right)$$

352 *Then with probability at least  $1 - \frac{\delta}{3}$ , every  $\ell$  satisfying  $|\ell| \leq \log(n/(2b_{\min}))$  that*  
 353 *MINIBATCHTOPDOWN $_{\mathcal{G}}(t, b, \mathbf{S})$  constructs has a minibatch,  $\mathbf{B} \sim \text{Batch}_b(\mathbf{S}, \ell)$ , of size at least*  
 354  *$b_{\min}$ .*

355 *Proof.* It suffices to show that the number of points in  $\mathbf{S}$  consistent with each of these  $\ell$  is at least  
 356  $b_{\min}$ . Fix any such  $\ell$  satisfying  $|\ell| \leq \log(n/(2b_{\min}))$ . The probability an element in  $\mathbf{S}$  is consistent  
 357 with  $\ell$  is at least  $\frac{2b_{\min}}{n}$ , meaning the expected number of points consistent with  $\ell$  is at least  $2b_{\min}$ . By  
 358 the multiplicative Chernoff bound,

$$\Pr \left[ \sum_{(x,y) \in \mathbf{S}} \mathbb{1}[x_i \text{ consistent with } \ell] < b_{\min} \right] \leq \exp_e \left( -\frac{1}{8} \cdot 2b_{\min} \right)$$

359 There are at most  $t$  leaves that MINIBATCHTOPDOWN $_{\mathcal{G}}(t, b, \mathbf{S})$  will ever estimate impurity gain for,  
 360 so as long as,

$$b_{\min} \geq 4 \cdot \log_e \left( \frac{3t}{\delta} \right),$$

361 with probability at least  $1 - \delta/3$ , all of them will have a minibatch of size at least  $b_{\min}$ .  $\square$

362 **Lemma B.2** (Batches are balanced). *With probability at least  $1 - \delta/3$ , there are at least  $\frac{b_{\min}}{4}$  points*  
 363  *$(x, y)$  in  $\mathbf{B}$  satisfying  $x_i = -1$  and  $\frac{b_{\min}}{4}$  points satisfying  $x_i = 1$ .*

364 *Proof.* The mini batch  $\mathbf{B}$  is formed by choosing at least  $b_{\min}$  points that are consistent with  $\ell$ , without  
 365 replacement, from  $\mathbf{S}$ , which is itself formed by taking points with replacement from  $\{\pm 1\}^d$ . This  
 366 means that the mini batch  $\mathbf{B}$  has at least  $b_{\min}$  points without replacement from  $\{\pm 1\}^d$ . Fix any  $\ell$  and  
 367 let  $b_{\text{true}}$  be the number of points in  $\mathbf{B}$ . By Hoeffding's inequality,

$$\Pr \left[ \left| \frac{b_{\text{true}}}{2} - (\text{Number of } (x, y) \in \mathbf{B} \text{ where } x_i = -1) \right| \geq \frac{b_{\text{true}}}{4} \right] \leq \exp_e \left( -\frac{b_{\text{true}}}{8} \right) \\ \leq \exp_e \left( -\frac{b_{\min}}{8} \right)$$

368 MINIBATCHTOPDOWN $_{\mathcal{G}}$  computes LocalGain $_{\mathcal{G}, \mathbf{B}}(\ell, i)$  for at most  $t$  different  $\ell$  and  $d$  different  $i$ ,  
 369 for a total of  $t \cdot d$  different computations. As long as

$$b_{\min} \geq 8 \cdot \log_e \left( \frac{3td}{\delta} \right),$$

370 then with probability at least  $1 - \delta/3$ , both  $\mathbf{B}[x_i = -1]$  and  $\mathbf{B}[x_i = 1]$  will have at least  $\frac{b_{\text{true}}}{4} \geq \frac{b_{\min}}{4}$   
 371 points.  $\square$

372 **Lemma B.3** (Batch size is logarithmic in  $td$ ). *For any  $f : \{\pm 1\}^d \rightarrow \{0, 1\}$  and  $n \in \mathbb{N}$ , let  $\mathbf{S}$  be a*  
 373 *size  $n$  sample of points  $(\mathbf{x}, f(\mathbf{x}))$  where  $\mathbf{x} \sim \{\pm 1\}^d$ . Furthermore, let  $\mathcal{G} : [0, 1] \rightarrow [0, 1]$  be any*  
 374  *$(C, \alpha)$ -Hölder continuous impurity function. For any  $\Delta > 0$ , and*

$$b \geq b_{\min} = \max \left( 8, 2 \cdot \left( \frac{2C}{\Delta} \right)^{\frac{2}{\alpha}} \right) \cdot \log_e \left( \frac{9td}{\delta} \right)$$

375 *with probability at least  $1 - \delta$ , any time MINIBATCHTOPDOWN $_{\mathcal{G}}(t, b, \mathbf{S})$  computes*  
 376 *LocalGain $_{\mathcal{G}, \mathbf{B}}(\ell, i)$  for  $|\ell| \leq \log(n/(2b_{\min}))$  for a mini batch  $\mathbf{B} \sim \text{Batch}_b(\mathbf{S}, \ell)$*

$$|\text{LocalGain}_{\mathcal{G}, \mathbf{B}}(\ell, i) - \text{LocalGain}_{\mathcal{G}, f}(\ell, i)| \leq \Delta$$

377 *Proof.* For any particular  $\ell, i$ , in order to compute  $\text{LocalGain}_{\mathcal{G}, \mathcal{B}}$  we need to estimate three expecta-  
 378 tions:

$$\begin{aligned} & \mathcal{G}(\mathbb{E}[f(\mathbf{x})]) \\ & \mathcal{G}(\mathbb{E}[f(\mathbf{x}) \mid \mathbf{x} \text{ reaches } \ell, \mathbf{x}_i = -1]) \\ & \mathcal{G}(\mathbb{E}[f(\mathbf{x}) \mid \mathbf{x} \text{ reaches } \ell, \mathbf{x}_i = 1]) \end{aligned}$$

379 Define  $\varepsilon_1, \varepsilon_2, \varepsilon_3$  to be the errors made in computing these expectations so that

$$\begin{aligned} \text{LocalGain}_{\mathcal{G}, \mathcal{B}}(\ell, i) & := \mathcal{G}(\mathbb{E}[f(\mathbf{x}) \mid \mathbf{x} \text{ reaches } \ell]) + \varepsilon_1 \\ & - \left(\frac{1}{2} \cdot \mathcal{G}(\mathbb{E}[f(\mathbf{x}) \mid \mathbf{x} \text{ reaches } \ell, \mathbf{x}_i = -1]) + \varepsilon_2\right) \\ & + \frac{1}{2} \cdot \mathcal{G}(\mathbb{E}[f(\mathbf{x}) \mid \mathbf{x} \text{ reaches } \ell, \mathbf{x}_i = 1]) + \varepsilon_3). \end{aligned}$$

380 Suppose that  $\varepsilon_1, \varepsilon_2, \varepsilon_3$  are each bounded as

$$|\varepsilon_j| \leq \left(\frac{\Delta}{2C}\right)^{\frac{1}{\alpha}} \quad (1)$$

381 Then, by the definition of Hölder continuous and triangle inequality,

$$\begin{aligned} & |\text{LocalGain}_{\mathcal{G}, \mathcal{B}}(\ell, i) - \text{LocalGain}_{\mathcal{G}, f}(\ell, i)| \\ & \leq |\mathcal{G}(\mathbb{E}[f(\mathbf{x}) \mid \mathbf{x} \text{ reaches } \ell]) + \varepsilon_1 - \mathcal{G}(\mathbb{E}[f(\mathbf{x}) \mid \mathbf{x} \text{ reaches } \ell])| \\ & \quad + \frac{1}{2} |\mathcal{G}(\mathbb{E}[f(\mathbf{x}) \mid \mathbf{x} \text{ reaches } \ell, \mathbf{x}_i = -1]) + \varepsilon_2 - \mathcal{G}(\mathbb{E}[f(\mathbf{x}) \mid \mathbf{x} \text{ reaches } \ell, \mathbf{x}_i = -1])| \\ & \quad + \frac{1}{2} |\mathcal{G}(\mathbb{E}[f(\mathbf{x}) \mid \mathbf{x} \text{ reaches } \ell, \mathbf{x}_i = 1]) + \varepsilon_3 - \mathcal{G}(\mathbb{E}[f(\mathbf{x}) \mid \mathbf{x} \text{ reaches } \ell, \mathbf{x}_i = 1])| \\ & \leq C \cdot (|\varepsilon_1|^\alpha + \frac{1}{2} \cdot |\varepsilon_2|^\alpha + \frac{1}{2} \cdot |\varepsilon_3|^\alpha) \\ & \leq \Delta \end{aligned}$$

382 Therefore, it is enough to show that for all  $\ell, i$ , the corresponding  $\varepsilon_1, \varepsilon_2, \varepsilon_3$  satisfy Equation (1).

383 By Lemma B.1 and Lemma B.2, with high probability all of these expectations are over at least  $\frac{b_{\min}}{4}$   
 384 terms. Given the above is true, we can use Hoeffding's inequality to bound each  $\varepsilon_j$ ,

$$\Pr \left[ |\varepsilon_j| > \left(\frac{\Delta}{2C}\right)^{\frac{1}{\alpha}} \right] \leq \exp_e \left( -2 \cdot \frac{b_{\min}}{4} \cdot \left(\frac{\Delta}{2C}\right)^{\frac{2}{\alpha}} \right).$$

385 There are a total of at most  $3td$  such  $\varepsilon_j$  we wish to bound. Setting  $b_{\min}$  to at least

$$2 \cdot \left(\frac{2C}{\Delta}\right)^{\frac{2}{\alpha}} \cdot \log_e \left(\frac{9td}{\delta}\right)$$

386 means all are bounded as desired with probability at least  $1 - \delta/3$ .  $\square$

## 387 B.2 Properties of MiniBatchTopDown

388 As we discussed in Section 2, a key component of [BLT20a]'s analysis is a proof that if error  $f(T_S^\circ) >$   
 389  $\text{opt}_s + \varepsilon$ , there must exist a leaf  $\ell^* \in T^\circ$  and a coordinate  $i^* \in [d]$  such that

$$\text{PurityGain}_{\mathcal{G}, f}(\ell^*, i^*) > \frac{\kappa \varepsilon^2}{32j(\log s)^2}. \quad (2)$$

390 Based on how we set  $\Delta$  in Lemma B.3,  $\text{MINIBATCHTOPDOWN}_{\mathcal{G}}$  will be able to estimate all local  
 391 gains to additive accuracy  $\pm O\left(\frac{\kappa \varepsilon^2}{\log(s)^2}\right)$ . That accuracy, in conjunction with just Equation (2), is *not*  
 392 sufficient to prove that  $\text{MINIBATCHTOPDOWN}_{\mathcal{G}}$  will produce a low error tree. Instead, we need the  
 393 following additional fact that [BLT20a] proved one step prior to showing Equation (2); in fact, it  
 394 implies Equation (2) but is stronger, and that strength is needed for our purposes.

395 **Fact B.4** (Showned during the proof of Theorem 2 of [BLT20a]). *Let  $T^\circ$  be any partial tree. For any*  
 396  *$f : \{\pm 1\}^d \rightarrow \{0, 1\}$  and  $\kappa$ -strongly concave impurity function  $\mathcal{G} : [0, 1] \rightarrow [0, 1]$ , if  $\text{error}_f(T_S^\circ) >$*   
 397  *$\text{opt}_s + \varepsilon$ , then*

$$\sum_{\text{leaves } \ell \in T^\circ} \max_{i \in [d]} (\text{PurityGain}_{\mathcal{G},f}(\ell, i)) > \frac{\kappa}{32} \cdot \left( \frac{\varepsilon}{\log s} \right)^2.$$

398 Fact B.4 implies Equation (2) because, if  $T^\circ$  is size  $j$  and the total purity gains of all of its leaves  
 399 is some value  $z$ , then at least one leaf has purity gain  $\frac{z}{j}$ . We use Fact B.4 to show that, whenever  
 400 `MINIBATCHTOPDOWN $\mathcal{G}$`  picks a leaf that is neither too deep nor too high in the tree, it has picked a  
 401 leaf and index with relatively large purity gain.

402 **Lemma B.5** (Medium depth splits are good.). *Choose any max depth  $D \in \mathbb{N}$ . Let  $f : \{\pm 1\}^d \rightarrow$*   
 403  *$\{0, 1\}$  be a monotone target function and  $\mathcal{G}$  be any  $\kappa$ -strongly concave and  $(C, \alpha)$ -Hölder continuous*  
 404 *impurity function. For any  $s \in \mathbb{N}$ ,  $\varepsilon, \delta \in (0, \frac{1}{2})$ , let  $t = s^{O(\log(s))/\varepsilon^2}$ , and  $\mathbf{S}$  be a set of  $n$  labeled*  
 405 *examples  $(\mathbf{x}, f(\mathbf{x}))$  where  $\mathbf{x} \sim \{\pm 1\}^d$  is uniform random,*

$$n = \Omega \left( \left( \frac{C^2 \log(s)^4}{\kappa^2 \varepsilon^4} \right)^{\frac{1}{\alpha}} \cdot \log \left( \frac{td}{\delta} \right) \cdot 2^D \right)$$

406 and

$$b = \Omega \left( \left( \frac{C^2 \log(s)^4}{\kappa^2 \varepsilon^4} \right)^{\frac{1}{\alpha}} \cdot \log \left( \frac{td}{\delta} \right) \right).$$

407 Then, with probability at least  $1 - \delta$ , the following holds for all iterations of  
 408 `MINIBATCHTOPDOWN $\mathcal{G}$` ( $t, b, \mathbf{S}$ ). If, at iteration  $j$ ,  $T^\circ$  satisfies,

$$\text{error}_f(T_{\text{Batch}_b(S)}^\circ) \geq \text{opt}_s + 2\varepsilon,$$

409 let  $(\ell^*, i^*)$  be the leaf and coordinate chosen to maximize the  $\text{PurityGain}_{\mathcal{G},\mathbf{B}}$ . Then, if

$$\log(j) - 2 \leq |\ell^*| \leq D,$$

410 then

$$\text{PurityGain}_{\mathcal{G},f}(\ell^*, i^*) > \frac{\kappa}{64} \cdot \frac{\varepsilon^2}{j(\log s)^2}.$$

411 *Proof.* For the values of  $n$  and  $b$  given in this lemma statement, using Lemma B.3, we have for  
 412  $\Delta = \frac{\kappa}{32 \cdot 10} \cdot \left( \frac{\varepsilon}{\log s} \right)^2$ , for all leaves with  $|\ell| \leq D$

$$|\text{LocalGain}_{\mathcal{G},\mathbf{B}}(\ell, i) - \text{LocalGain}_{\mathcal{G},f}(\ell, i)| \leq \frac{\kappa}{32 \cdot 10} \cdot \left( \frac{\varepsilon}{\log s} \right)^2 \quad (3)$$

413 with probability atleast  $1 - \delta$ .

414 Since  $\text{error}_f((T^\circ)_{\text{Batch}_b(S)}) \geq \text{opt}_s + 2\varepsilon$  and using lemma B.1, we know that  $\text{error}_f((T^\circ)_S) \geq$   
 415  $\text{opt}_s + \varepsilon$  since the batch size is large enough, we can use Fact B.4 to lower bound the estimated purity  
 416 gain of  $\ell^*$  and  $i^*$ . Let  $c = \frac{\kappa}{32}$

$$\begin{aligned} \sum_{\text{leaves } \ell \in T^\circ} \max_{i \in [d]} (\text{PurityGain}_{\mathcal{G},f}(\ell, i)) &> c \cdot \left( \frac{\varepsilon}{\log s} \right)^2 && \text{Fact B.4} \\ \sum_{\text{leaves } \ell \in T^\circ} 2^{-|\ell|} \cdot \max_{i \in [d]} (\text{LocalGain}_{\mathcal{G},f}(\ell, i)) &> c \cdot \left( \frac{\varepsilon}{\log s} \right)^2 \\ \sum_{\text{leaves } \ell \in T^\circ} 2^{-|\ell|} \cdot \max_{i \in [d]} (\text{LocalGain}_{\mathcal{G},\mathbf{B}}(\ell, i)) &> \frac{9c}{10} \cdot \left( \frac{\varepsilon}{\log s} \right)^2 && \text{Equation (3) and } \sum_{\ell} 2^{-|\ell|} = 1 \\ \sum_{\text{leaves } \ell \in T^\circ} \max_{i \in [d]} (\text{PurityGain}_{\mathcal{G},\mathbf{B}}(\ell, i)) &> \frac{9c}{10} \cdot \left( \frac{\varepsilon}{\log s} \right)^2 \end{aligned}$$

417 Since there are  $j$  leaves in  $T^\circ$  and  $\ell^*, i^*$  are chosen to maximize  $\text{PurityGain}_{\mathcal{G}, \mathcal{B}}(\ell^*, i^*)$ ,

$$\text{PurityGain}_{\mathcal{G}, \mathcal{B}}(\ell^*, i^*) > \frac{9c}{10j} \cdot \left( \frac{\varepsilon}{\log s} \right)^2.$$

418 Next, we show that since  $\ell^*$  is sufficiently far down in the tree, then the estimated purity gain and  
419 true purity gain are close.

$$\begin{aligned} & |\text{PurityGain}_{\mathcal{G}, \mathcal{B}}(\ell^*, i^*) - \text{PurityGain}_{\mathcal{G}, f}(\ell^*, i^*)| \\ &= 2^{-|\ell^*|} \cdot |\text{LocalGain}_{\mathcal{G}, \mathcal{B}}(\ell^*, i^*) - \text{LocalGain}_{\mathcal{G}, f}(\ell^*, i^*)| \\ &\leq 2^{-|\ell^*|} \cdot \frac{c}{10} \cdot \left( \frac{\varepsilon}{\log s} \right)^2 && \text{eq. (3)} \\ &\leq \frac{4}{j} \cdot \frac{c}{10} \cdot \left( \frac{\varepsilon}{\log s} \right)^2 && |\ell^*| \geq \log(j) - 2 \end{aligned}$$

420 By triangle inequality, we have that  $\text{PurityGain}_{\mathcal{G}, f}(\ell^*, i^*) > \frac{c}{2j} \cdot \left( \frac{\varepsilon}{\log s} \right)^2$ , the desired result.

421 □

422 Given that we are only guaranteed to make good progress on splits that are neither too deep nor too  
423 shallow, we will need to deal with both possibilities. First, we show that if we ever wanted to make  
424 too deep a split, we would already be done.

425 **Lemma B.6** (Can stop at very large depth.). *Let  $f : \{\pm 1\}^d \rightarrow \{0, 1\}$  be a monotone target function  
426 and  $\mathcal{G}$  be any  $\kappa$ -strongly concave and  $(C, \alpha)$ -Hölder continuous impurity function. For any  $s \in \mathbb{N}$ ,  
427  $\varepsilon, \delta \in (0, \frac{1}{2})$ , let*

$$t = s^{\Theta(\log(s))/(\kappa\varepsilon^2)}, \quad (4)$$

428 set the max depth to

$$D = \lfloor \log(t) + \log \log t \rfloor, \quad (5)$$

429 let  $\mathcal{S}$  be a set of  $n$  labeled examples  $(\mathbf{x}, f(\mathbf{x}))$  where  $\mathbf{x} \sim \{\pm 1\}^d$  is uniform random,

$$n = \Omega \left( \left( \frac{C^2 \log(s)^4}{\kappa^2 \varepsilon^4} \right)^{\frac{1}{\alpha}} \log \left( \frac{td}{\delta} \right) \cdot 2^D \right) = \text{poly}_{\alpha, \kappa, C}(t, \log(d), \log(1/\delta)),$$

430 and batch size at least

$$b = \Omega \left( \left( \frac{C^2 \log(s)^4}{\kappa^2 \varepsilon^4} \right)^{\frac{1}{\alpha}} \log \left( \frac{td}{\delta} \right) \right).$$

431 Let  $T_1^\circ, T_2^\circ, \dots, T_t^\circ$  be the size  $1, 2, \dots, t$  partials trees that  $\text{MINIBATCHTOPDOWN}_{\mathcal{G}}(t, b, \mathcal{S})$  builds.  
432 With probability  $1 - \delta$  over the randomness of  $\mathcal{S}$  and the random batches, for any  $k \in [t]$ , if  $T_k^\circ$  has  
433 depth more than  $D$ , then

$$\text{error}_f((T_k^\circ)_{\text{Batch}_b(\mathcal{S})}) \leq \text{opt}_s + 2\varepsilon. \quad (6)$$

434 *Proof.* Let  $k$  be chosen so that  $T_k^\circ$  has depth more than  $D$ . For some  $j \leq k$ , there was a leaf  $\ell^* \in T_j^\circ$   
435 that was split, satisfying,

$$\begin{aligned} |\ell^*| &= \lfloor \log(t) + \log \log t \rfloor - 1 \\ &= \lfloor \log(t) + \log(\Theta((\log s)^2/(\kappa\varepsilon^2))) \rfloor - 1. \end{aligned}$$

436 For any  $i \in [d]$ ,

$$\begin{aligned} \text{PurityGain}_{\mathcal{G}, f}(\ell^*, i) &= 2^{-|\ell^*|} \text{LocalGain}_{\mathcal{G}, f}(\ell^*, i) \\ &\leq 2^{-|\ell^*|} \\ &\leq \frac{1}{t} \cdot \Theta \left( \left( \frac{(\log s)^2}{\kappa\varepsilon^2} \right)^{-1} \right) \\ &\leq \frac{1}{j} \cdot \Theta \left( \frac{\kappa\varepsilon^2}{(\log s)^2} \right). \end{aligned} \quad (7)$$

437 Note that the constant in Equation (7) is inversely related to the constant in the exponent of Equa-  
 438 tion (4). In Lemma B.5, we showed that if  $\text{error}_f((T_j^\circ)_{\text{Batch}_b(S)}) \geq \text{opt}_s + \varepsilon$ , then for some  
 439  $i^* \in [d]$ ,

$$\text{PurityGain}_{\mathcal{G},f}(\ell^*, i^*) = \Omega\left(\frac{\kappa\varepsilon^2}{j(\log s)^2}\right). \quad (8)$$

440 If we choose the constant in Equation (7) sufficiently low, which can be done by making the constant in  
 441 Equation (4) sufficiently high, then that equation can not be satisfied at the same time as Equation (8).  
 442 Therefore, it must be that  $\text{error}_f((T_j^\circ)_{\text{Batch}_b(S)}) < \text{opt}_s + \varepsilon$ . Since  $j < k$ , and adding splits can only  
 443 increase error by at most  $\varepsilon$ , it must also be the case that  $\text{error}_f((T_k^\circ)_{\text{Batch}_b(S)}) < \text{opt}_s + 2\varepsilon$ .  $\square$

444 We next show that before Lemma B.6 kicks in, most splits are sufficiently deep to make good progress.

445 **Lemma B.7** (Few splits are shallow). *Let  $k = 2^a$  be any power of 2 and  $T_1^\circ, \dots, T_k^\circ$  be a series of*  
 446 *bare trees of size  $1, \dots, k$  respectively where  $T_{j+1}$  is formed by splitting  $\ell_j \in T_j$ . Then,*

$$\sum_{j=1}^k \mathbb{1}[|\ell_j| < \log(j) - 2] \leq \frac{k}{4}$$

447 *Proof.* First, since for all  $j = 1, \dots, k, j \leq k$ , we can bound,

$$\sum_{j=1}^k \mathbb{1}[|\ell_j| < \log(j) - 2] \leq \sum_{j=1}^k \mathbb{1}[|\ell_j| < \log(k) - 2] = \sum_{j=1}^k \mathbb{1}[|\ell_j| < a - 2].$$

448 If  $\ell_j$ , a leaf of  $T_j^\circ$ , has depth less than  $a - 2$ , then it is also an internal node of  $T_k^\circ$  with depth less  
 449 than  $a - 2$ . There are at most  $2^{a-2} - 1$  nodes in any tree of depth less than  $a - 2$ . Therefore,

$$\sum_{j=1}^k \mathbb{1}[|\ell_j| \leq a - 2] \leq 2^{a-2} - 1 \leq \frac{k}{4} - 1 \leq \frac{k}{4}.$$

450

$\square$

### 451 B.3 Final proof of Theorem 1

452 *Proof.* MINIBATCHTOPDOWN $\mathcal{G}$  builds a series of bare trees,  $T_1^\circ, T_2^\circ, \dots, T_t^\circ$ , where  $T_j$  has size  $j$ .  
 453 We wish to prove that  $\text{error}_f((T_t^\circ)_{\text{Batch}_b(S)}) \leq \text{opt}_s + 3\varepsilon$  (In the end, we can choose  $\varepsilon$  appropriately  
 454 to get error  $\text{opt}_s + \varepsilon$ ). To do so, we consider two cases.

455 **Case 1:** There is some  $k < t$  for which  $\text{error}_f((T_k^\circ)_{\text{Batch}_b(S)}) \leq \text{opt}_s + 2\varepsilon$ .  
 456 Since splitting more variables of  $T_k$  can only increase its error by at most  $\varepsilon$ ,

$$\text{error}_f((T_t^\circ)_{\text{Batch}_b(S)}) \leq \text{error}_f((T_k^\circ)_{\text{Batch}_b(S)}) + \varepsilon \leq \text{opt}_s + 3\varepsilon,$$

457 which is the desired result.

458 **Case 2:** There is no  $k < t$  for which  $\text{error}_f((T_k^\circ)_{\text{Batch}_b(S)}) \leq \text{opt}_s + \varepsilon$ .

459

460 In this case, we use Lemma B.5 to ensure we make good progress. Lemma B.5 only applies when the  
 461 tree has depth at most  $D = \log t + \log \log t$ . Luckily, Lemma B.6 ensures that if the tree has depth  
 462 more than  $D$ , then we are ensured that  $\text{error}_f((T_t^\circ)_{\text{Batch}_b(S)}) \leq \text{opt}_s + 2\varepsilon$ , and so are done. For the  
 463 remainder of this proof, we assume all partial trees have depth at most  $D$ .

464 We will show that  $\mathcal{G}$ -impurity( $T_t^\circ$ ) = 0, which means that  $\text{error}_f((T_t^\circ)_{\text{Batch}_b(S)}) = 0 \leq \text{opt}_s + 2\varepsilon$ ,  
 465 also proving the desired result. For  $j = 1, \dots, t - 1$ , let  $\ell_j$  be the leaf of  $T_j$  that is split, and  $i_j$  be the  
 466 coordinate placed at  $\ell_j$  to form  $T_{j+1}$ . Then,

$$\mathcal{G}\text{-impurity}(T_t^\circ) = \mathcal{G}\text{-impurity}(T_1^\circ) - \sum_{j=1}^{t-1} \text{PurityGain}_{\mathcal{G},f}(\ell_j, i_j).$$

467 Since  $\mathcal{G}$ -impurity( $T_1^\circ$ )  $\leq 1$  and our goal is to show that  $\mathcal{G}$ -impurity( $T_{t+1}^\circ$ ) = 0, it is sufficient to  
 468 show that  $\sum_{j=1}^t \text{PurityGain}_{\mathcal{G},f}(\ell_j, i_j) \geq 1$ . Lemma B.5 combined with Fact B.4,

$$\sum_{j=1}^t \text{PurityGain}_{\mathcal{G},f}(\ell_j, i_j) \geq \sum_{j=1}^t \mathbf{1}[|\ell_j| \geq \log(j) - 2] \cdot \frac{\kappa}{64j} \cdot \left(\frac{\varepsilon}{\log s}\right)^2$$

469 We break the above summation into chunks from  $j = (2^a + 1)$  to  $j = 2^{a+1}$ , integer  $a \leq \log(t)$ . In  
 470 such a chunk, there are  $2^a$  choices for  $j$ . By Lemma B.7, we know that for at most  $2^{a+1}/4 = 2^a/2$   
 471 of those  $j$  is  $\mathbf{1}[|\ell_j| < \log(j) - 2]$ . Therefore,

$$\begin{aligned} \sum_{j=2^{a+1}}^{2^{a+1}} \mathbf{1}[|\ell_j| \leq \log(j) - 2] \cdot \frac{\kappa}{64j} \cdot \left(\frac{\varepsilon}{\log s}\right)^2 &\geq \frac{2^a}{2} \cdot \frac{\kappa}{64 \cdot (2^{a+1})} \cdot \left(\frac{\varepsilon}{\log s}\right)^2 \\ &= \frac{\kappa}{256} \cdot \left(\frac{\varepsilon}{\log s}\right)^2 \end{aligned}$$

472 Summing up  $\frac{256}{\kappa} \cdot \left(\frac{\log s}{\varepsilon}\right)^2$  such chunks gives a sum of at least 1. Therefore, for

$$t = \exp\left(\Omega\left(\frac{(\log s)^2}{\kappa \varepsilon^2}\right)\right)$$

473 it must be the case that  $\mathcal{G}$ -impurity( $T_{t+1}^\circ$ ) = 0, proving the desired result.  $\square$

## 474 C Estimating the size of a decision tree

475 In this section, we design a decision tree size estimator. This size estimator only needs to inspect a  
 476 small number of random strands from the decision tree. It is unbiased, and as long as the decision  
 477 tree has a bounded max depth, obeys concentration bounds shown in Lemma C.1.

478 **Lemma C.1** (Size estimator). *For any  $\Delta, \delta > 0$  and size- $s$  decision tree  $T$ , let  $\ell^*$  be the deepest leaf  
 479 in  $T$  and*

$$m = \frac{(2^{|\ell^*|})^2}{2\Delta^2} \cdot \ln\left(\frac{2}{\delta}\right).$$

480 *Choose  $\mathbf{x}_1, \dots, \mathbf{x}_m$  uniformly random from  $\{\pm 1\}^d$  and define the estimator*

$$e := \frac{1}{m} \sum_{i=1}^m 2^{|\ell_T(\mathbf{x}_i)|}.$$

481 *With probability at least  $1 - \delta$ ,*

$$|e - s| \leq \Delta.$$

482 *Proof.* We first show that  $\mathbb{E}[e] = s$ .

$$\begin{aligned} \mathbb{E}[e] &= \mathbb{E}_{\mathbf{x} \sim \{\pm 1\}^d} \left[ 2^{|\ell_T(\mathbf{x})|} \right] \\ &= \sum_{\text{leaves } \ell \in T} \Pr[\mathbf{x} \text{ reaches } \ell] \cdot 2^{|\ell|} \\ &= \sum_{\text{leaves } \ell \in T} \frac{1}{2^{|\ell|}} \cdot 2^{|\ell|} \\ &= s, \end{aligned}$$

483 where the last equality is due to the fact that a size- $s$  tree has  $s$  leaves. Furthermore,  $e$  is the sum of  
 484  $m$  independent random variables bounded between 0 and  $2^{|\ell^*|}$ . Therefore, we can apply Hoeffding's  
 485 inequality,

$$\Pr[|e - s| \geq \Delta] \leq 2 \exp_e\left(-\frac{2m\Delta^2}{(2^{|\ell^*|})^2}\right).$$

486 Plugging in  $m$  proves the desired result.  $\square$



487 **D Provable guarantees for LOCALLEARNER**

488 In order to facilitate comparisons between the output of LOCALLEARNER $\mathcal{G}$  and  
 489 MINIBATCHTOPDOWN $\mathcal{G}$ , we will define another algorithm, TOPDOWNSIZEESTIMATE $\mathcal{G}$ , that  
 490 shares some elements with LOCALLEARNER $\mathcal{G}$  and some elements with MINIBATCHTOPDOWN $\mathcal{G}$ .

TOPDOWNSIZEESTIMATE $\mathcal{G}(t, b, S)$ :

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

Define  $D := \log t + \log \log t$ .

Let  $B_{\text{strands}}^\circ$  be  $b$  uniform random points from  $\{\pm 1\}^d$ .

Initialize  $e := 1$ , our size estimate.

while ( $e < t$ ) {

1. *Score*: For each leaf  $\ell \in T^\circ$  of depth at most  $D$ , draw  $B \sim \text{Batch}_b(S, \ell)$ . For each coordinate  $i \in [d]$ , compute:
 
$$\text{PurityGain}_{\mathcal{G}, B}(\ell, i) := 2^{-|\ell|} \cdot \text{LocalGain}_{\mathcal{G}, B}(\ell, i), \text{ where}$$

$$\text{LocalGain}_{\mathcal{G}, B}(\ell, i) := \mathcal{G}(\mathbb{E}[f(\mathbf{x})]) - \left( \frac{1}{2} \mathcal{G}(\mathbb{E}[f(\mathbf{x}) \mid \mathbf{x}_i = -1]) + \frac{1}{2} \mathcal{G}(\mathbb{E}[f(\mathbf{x}) \mid \mathbf{x}_i = 1]) \right),$$
 where the expectations are with respect to  $(\mathbf{x}, f(\mathbf{x})) \sim B$ .
2. *Split*: Let  $(\ell^*, i^*)$  be the tuple that maximizes  $\text{PurityGain}_{\mathcal{G}, B}(\ell, i)$ . Grow  $T^\circ$  by splitting  $\ell^*$  with a query to  $x_{i^*}$ .
3. *Estimate size*: Update our size estimate to
 
$$e = \mathbb{E}_{\mathbf{x} \in X} [2^{|\ell_{T^\circ}(\mathbf{x})|}]$$

}

For each leaf  $\ell \in T^\circ$ , draw  $B \sim \text{Batch}_b(S, \ell)$  and label  $\ell$  with  $\text{round}(\mathbb{E}_{(\mathbf{x}, f(\mathbf{x})) \sim B} [f(\mathbf{x})])$ .

Figure 4: TOPDOWNSIZEESTIMATE $\mathcal{G}$  takes as input a size parameter  $t$ , a minibatch size  $b$ , and a labeled dataset  $S$ . It outputs a size- $t'$  decision tree hypothesis for  $f$ , where  $t'$  is close to  $t$ .

491 **Comparison between MINIBATCHTOPDOWN $\mathcal{G}$  and TOPDOWNSIZEESTIMATE $\mathcal{G}$ :** The only  
 492 difference between MINIBATCHTOPDOWN $\mathcal{G}$  and TOPDOWNSIZEESTIMATE $\mathcal{G}$  is the stopping cri-  
 493 terion. MINIBATCHTOPDOWN $\mathcal{G}$  stops when the size of  $T^\circ$  is exactly  $t$ . On the other hand,  
 494 TOPDOWNSIZEESTIMATE $\mathcal{G}$  estimates the size of  $T^\circ$  using the estimator from Appendix C and  
 495 stops when this size estimate is at least  $t$ .

496 **Comparison between LOCALLEARNER $\mathcal{G}$  and MINIBATCHTOPDOWN $\mathcal{G}$ :** For any  $t, S, b, x^*$   
 497 that are valid inputs to LOCALLEARNER $\mathcal{G}$ , we compare the following two procedures.

- 498 1. Running TOPDOWNSIZEESTIMATE $\mathcal{G}(t, b, S)$  to get a decision tree,  $T$ , and then computing  
 499  $T(x^*)$ .
- 500 2. Only running LOCALLEARNER $\mathcal{G}(t, b, S, x^*)$ .

501 We claim the output from the above two procedures is identical (given Footnote 2).  
 502 TOPDOWNSIZEESTIMATE $\mathcal{G}$  expands all paths in the tree its building, whereas LOCALLEARNER $\mathcal{G}$   
 503 only expands paths that are pertinent to either the input  $x^*$ , or inputs in  $B_{\text{strands}}^\circ$ , which are used to  
 504 compute the size estimate. Aside from that, both of the above procedures are identical. Furthermore,  
 505 paths not containing  $x^*$  nor any inputs in  $B_{\text{strands}}^\circ$  have no effect on how the tree eventually labels  $x^*$ .  
 506 Therefore, the output of the two above procedures is identical, though LOCALLEARNER $\mathcal{G}$  is more  
 507 efficient as it only computes necessary paths.

508 Combining the above observations, we are able to prove the formal version of Theorem 2.

509 **Theorem 5** (Formal version of Theorem 2). *Let  $f : \{\pm 1\}^d \rightarrow \{0, 1\}$  be a target function,  $\mathcal{G}$  be an*  
 510 *impurity function, and  $S^\circ$  be an unlabeled training set.*

511 *For all  $t \in \mathbb{N}$  and  $\eta, \delta \in (0, \frac{1}{2})$ , if the minibatch size is at least*

$$b = \Omega\left(\frac{(\log t)^2}{\eta^2} \cdot \log\left(\frac{t}{\delta}\right)\right),$$

512 *then with probability at least  $1 - \delta$  over the randomness of  $\mathbf{B}_{\text{strands}}^\circ$ , there is some  $t' \in [t - \eta t, t + \eta t]$*   
 513 *for which the following holds. For all  $x^* \in \{\pm 1\}^d$ ,  $\text{LOCALLEARNER}_{\mathcal{G}}(t, b, S^\circ, x^*)$  labels*

$$q = O(b^2 \log t)$$

514 *points within  $S^\circ$  and returns  $T(x^*)$ , where  $T$  is the size- $t'$  decision tree hypothesis that*  
 515  *$\text{MINIBATCHTOPDOWN}_{\mathcal{G}}(t', b, S)$  would construct, and  $S$  is the labeled dataset obtained by la-*  
 516 *beling all of  $S^\circ$  with  $f$ 's values.*

517 We break the proof of Theorem 5 into two pieces. First, we show that it labels only  $O(b^2 \log t)$  points  
 518 within  $S^\circ$ , and then the rest.

519 **Lemma D.1** (Label efficiency of  $\text{LOCALLEARNER}_{\mathcal{G}}$ ). *Let  $f : \{\pm 1\}^d \rightarrow \{0, 1\}$  be a target function,*  
 520  *$\mathcal{G}$  be an impurity function, and  $S^\circ$  be an unlabeled training set. For any  $b, t \in \mathbb{N}$  and  $x^* \in \{\pm 1\}^d$ ,*  
 521  *$\text{LOCALLEARNER}_{\mathcal{G}}(t, b, S^\circ, x^*)$  labels at most*

$$q = O(b^2 \log t)$$

522 *points within  $S^\circ$ .*

523 *Proof.* It is sufficient for us to show that  $\text{LOCALLEARNER}_{\mathcal{G}}$  labels at most  $O(b \log t)$  batches.  
 524  $\text{LOCALLEARNER}_{\mathcal{G}}$  builds a series of bare trees  $T_1^\circ, \dots, T_{t'}^\circ$ . During the while loop, the number of  
 525 batches it labels is equal to nodes in the following set

$$L := \bigcup_{j=1}^{t'} \left\{ \ell_{T_j^\circ}(x) : x \in \mathbf{B}_{\text{strands}}^\circ \cup \{x^*\}, |\ell_{T_j^\circ}(x)| \leq D \right\}$$

526 Consider a single  $x \in \mathbf{B}_{\text{strands}}^\circ \cup \{x^*\}$ , and define

$$L(x) := \left\{ \ell_{T_j^\circ}(x) : j \in [t'], |\ell_{T_j^\circ}(x)| \leq D \right\}.$$

527 Every node in  $L(x)$  has depth at most  $D$ , and there is at most one node in  $L(x)$  per depth. Therefore,  
 528  $|L(x)| \leq D$ , and

$$\begin{aligned} |L| &\leq \sum_{x \in \mathbf{B}_{\text{strands}}^\circ \cup \{x^*\}} |L(x)| \\ &\leq (b+1)D \\ &= O(b \log t). \end{aligned}$$

529 Therefore,  $\text{LOCALLEARNER}_{\mathcal{G}}$  labels only  $O(b \log t)$  batches during the while loop. After the while  
 530 loop, it labels at most 1 additional batches. Therefore, it labels a total of  $O(b \log t)$  batches which  
 531 requires labeling  $O(b^2 \log t)$  points.  $\square$

532 We next prove the remainder of Theorem 5.

533 *Proof.* Let  $T$  be the tree that  $\text{TOPDOWNSIZEESTIMATE}_{\mathcal{G}}(t, b, S)$  produces. In the comparison  
 534 between  $\text{LOCALLEARNER}_{\mathcal{G}}$  and  $\text{TOPDOWNSIZEESTIMATE}_{\mathcal{G}}$ , we established that, for all  $x^* \in$   
 535  $\{\pm 1\}^d$ ,

$$\text{LOCALLEARNER}_{\mathcal{G}}(t, b, S, x^*) = T(x^*).$$

536 Set  $t' = |T|$ . Then,  $T$  is also the output of  $\text{MINIBATCHTOPDOWN}_{\mathcal{G}}(t', b, S)$ , as desired. Next, we  
 537 prove that  $t' \in [t - \eta t, t + \eta t]$  with probability at least  $1 - \delta$ .

538 Let  $T_1^\circ, T_2^\circ, \dots, T_{t'}^\circ$  be the bare trees of size  $1, 2, \dots, t'$  that  $\text{TOPDOWNSIZEESTIMATE}_{\mathcal{G}}(t, b, \mathbf{S})$  pro-  
 539 duces, and let  $e_1, e_2, \dots, e_{t'}$  be the corresponding size estimates. Since  $\text{TOPDOWNSIZEESTIMATE}_{\mathcal{G}}$   
 540 halts when the size estimate is at least  $t$ ,

$$e_{t'} \geq t \quad \text{and} \quad e_{t'-1} < t.$$

541 We set  $\Delta := \eta t$  and wish, for all  $1 \leq j \leq t + \eta t$ , that  $e_j$  estimate the size of  $T_j^\circ$  to accuracy  $\pm \Delta$ .  
 542 Since the size of  $T_j^\circ$  is  $j$ , we equivalently wish for

$$|e_j - j| \leq \Delta \quad \text{for all } j = 1, \dots, t + \eta t. \quad (9)$$

543 Each  $T_j^\circ$  has max depth at most  $\log t + \log \log t$ . By Lemma C.1 and a union bound over all  $t + \eta t$   
 544 different  $j$ , we can guarantee that Equation (9) holds with probability at least  $1 - \delta$  if we set

$$\begin{aligned} b &\geq \frac{(2^{\log t + \log \log t})^2}{2\Delta^2} \cdot \ln \left( \frac{2t(1 + \eta)}{\delta} \right) \\ &= \Omega \left( \frac{(t \log t)^2}{(\eta t)^2} \cdot \log \left( \frac{t}{\delta} \right) \right) \\ &= \Omega \left( \frac{(\log t)^2}{\eta^2} \cdot \log \left( \frac{t}{\delta} \right) \right). \end{aligned}$$

545 Therefore, for the  $b$  we set in Theorem 5, Equation (9) holds with probability at least  $1 - \delta$ . For the  
 546 remainder of this proof, we suppose it holds and then show that the  $t' \in [t - \eta t, t + \eta t]$ . We first  
 547 show that  $t' \leq t + \eta t$ . By Equation (9), for  $j = t + \eta t$ ,

$$\begin{aligned} e_{(t+\eta t)} &\geq (t + \eta t) - \Delta \\ &\geq t. \end{aligned}$$

548 Recall that  $t'$  is the lowest integer such that  $e_{t'} \geq t$ . Therefore,  $t' \leq t + \eta t$ . We next show that  
 549  $t' \geq t - \eta t$ . By Equation (9) for  $j = t' \leq t + \eta t$ ,

$$\begin{aligned} t' &\geq e_{t'} - \Delta \\ &\geq t - \eta t. \end{aligned}$$

550 Therefore, Equation (9) implies  $t' \in [t - \eta t, t + \eta t]$  proving that with probability at least  $1 - \delta$ .

551 □

552 Finally, we show that the following algorithm estimates learnability.

EST $_{\mathcal{G}}(t, b, S^\circ, S_{\text{test}})$ :

Return  $\frac{1}{|S_{\text{test}}|} \sum_{(x,y) \in S_{\text{test}}} \mathbb{1}[\text{LOCALLEARNER}_{\mathcal{G}}(t, b, S^\circ, x) \neq y]$

Figure 5: EST $_{\mathcal{G}}$  takes as input a size parameter  $t$ , a minibatch size  $b$ , and an unlabeled dataset  $S^\circ$  and labeled test set  $S_{\text{test}}$ . It outputs the error of the tree returned by  $\text{MINIBATCHTOPDOWN}(t', b, S)$  with respect to  $S_{\text{test}}$ , where  $S$  is the labeled version of  $S^\circ$  and  $t'$  is close to  $t$ . As in Footnote 2, the random outcome of  $\mathbf{B}_{\text{strands}}^\circ$  and the minibatches should be consistent across all runs of  $\text{LOCALLEARNER}_{\mathcal{G}}$ .

553 **Theorem 6** (Formal version of Theorem 3). *Let  $f : \{\pm 1\}^d \rightarrow \{0, 1\}$  be a target function,  $\mathcal{G}$  be an*  
 554 *impurity function,  $S^\circ$  be an unlabeled training set, and  $S_{\text{test}}$  be a labeled test set.*

555 *For all  $t \in \mathbb{N}$  and  $\eta, \delta \in (0, \frac{1}{2})$ , if the minibatch size  $b$  is as in Theorem 5, then with probability at*  
 556 *least  $1 - \delta$  over the randomness of  $\mathbf{B}_{\text{strands}}^\circ$ ,  $\text{EST}_{\mathcal{G}}(t, b, S^\circ, S_{\text{test}})$  labels*

$$q = O(|S_{\text{test}}| \cdot b \log t + b^2 \log t)$$

557 *points within  $S^\circ$  and returns*

$$\text{error}_{S_{\text{test}}}(T) := \Pr_{(\mathbf{x}, \mathbf{y}) \sim S_{\text{test}}} [T(\mathbf{x}) \neq \mathbf{y}],$$

558 *where  $T$  is as in Theorem 5.*

559 *Proof.* Based on Theorem 5,  $\text{EST}_{\mathcal{G}}$  returns the desired result, so we only need to prove it labels  
 560 few points within  $S^\circ$ . As in Footnote 2, the same  $\mathbf{B}_{\text{strands}}^\circ$  are chosen across multiple runs of  
 561  $\text{LOCALLEARNER}_{\mathcal{G}}$ . As shown in Lemma D.1, the total number of points it labels is  $O(b \log t) *$   
 562  $m$ , where  $m$  is the number of strands built.  $\text{EST}_{\mathcal{G}}$  needs to build  $b$  strands for points within  
 563  $\mathbf{B}_{\text{strands}}^\circ$  and  $|S_{\text{test}}|$  strands for the points within  $S_{\text{test}}$ . As long as it caches its labels across runs of  
 564  $\text{LOCALLEARNER}_{\mathcal{G}}$ , the total labels used will be

$$q = O(b \log t) \cdot (b + |S_{\text{test}}|) = O(|S_{\text{test}}| \cdot b \log t + b^2 \log t).$$

565

□