

Optimal Sparse Decision Trees

Xiyang Hu
Duke University
Durham, NC 27708
xiyang.hu@duke.edu

Cynthia Rudin
Duke University
Durham, NC 27708
cynthia@cs.duke.edu

Margo Seltzer
The University of British Columbia
Vancouver, BC V6T 1Z4, Canada
mseltzer@cs.ubc.ca

ABSTRACT

Decision tree algorithms have been among the most popular algorithms for interpretable (transparent) machine learning since the early 1980's. The problem that has plagued decision tree algorithms since their inception is their lack of optimality, or lack of guarantees of closeness to optimality: decision tree algorithms are often greedy or myopic, and sometimes produce unquestionably suboptimal models. Hardness of decision tree optimization is both a theoretical and practical obstacle, and even careful mathematical programming approaches have not been able to solve these problems efficiently. This work introduces the first practical algorithm for optimal decision trees for binary variables. The algorithm is a co-design of analytical bounds that reduce the search space and modern systems techniques, including data structures and a custom bit-vector library. We highlight possible steps to improving the scalability and speed of future generations of this algorithm based on insights from our theory and experiments.

CCS CONCEPTS

• **Computing methodologies** → **Discrete space search; Classification and regression trees.**

KEYWORDS

Decision trees; Optimization; Interpretable models

ACM Reference Format:

Xiyang Hu, Cynthia Rudin, and Margo Seltzer. 2019. Optimal Sparse Decision Trees. In *ACM*, New York, NY, USA, 11 pages. <https://doi.org/00.0000/0000000.0000000>

1 INTRODUCTION

Interpretable machine learning has been growing in importance, as society has started to realize the dangers of using black box models for high stakes decisions: complications with confounding have haunted our medical machine learning models [22], bad predictions from black boxes have announced to millions of people that their dangerous levels of air pollution were safe [15], high-stakes credit risk decisions are being made without proper justification, and black box risk predictions have been wreaking havoc with the perception of fairness of our criminal justice system [10]. In all of these applications – medical imaging, pollution modeling, recidivism risk, credit scoring – interpretable models have been created that are just as accurate as the best black box models (by the CDC, Arnold Foundation, and others). However, such interpretable models are not generally easy to construct. If we want people to replace their

black box models with interpretable models, the tools to build these interpretable models must first exist.

Decision trees are one of the leading forms of interpretable models. Despite several attempts over the last several decades to improve the optimality of decision tree algorithms, the CART [7] and C4.5 [19] decision tree algorithms (and other greedy tree-growing variants) have remained as dominant methods in practice for decision trees. CART and C4.5 grow decision trees from the top down without backtracking, which means that if a suboptimal split was introduced near the top of the tree, the algorithm could spend many extra splits trying to undo the mistake it made at the top, leading to less-accurate and less-interpretable trees. Problems with greedy splitting and pruning have been known since the early 1990's, when mathematical programming tools had started to be used for creating optimal binary-split decision trees [3, 4], in a line of work [5, 6, 16, 18] until the present [20], however, these techniques use all-purpose tools, and tend not to scale to realistically-sized problems unless simplified to trees of a specific form. Other works [11] make overly strong assumptions (e.g., independence of all variables) to ensure optimal trees are produced using greedy algorithms.

Our work takes a different approach to optimal decision trees than mathematical programming, greedy methods, brute force, or strong assumptions. We use a composite of analytical bounds to perform massive pruning of the search space, specialized data structures to store intermediate computations and symmetries, and a bit-vector library to evaluate decision trees more quickly on the data, along with fast search policies, computational reuse and other systems-based techniques in order to achieve optimal trees according to a linear balance between accuracy and number of leaves. Despite the hardness of finding optimal solutions, our algorithm is able to locate optimal trees and prove optimality (or closeness of optimality) in reasonable amounts of time for datasets of the sizes used in the criminal justice system (tens of thousands or millions of observations, tens of features). Our algorithm is the first practical optimal decision tree algorithm to achieve solutions for problems of these sizes.

Because we find provably optimal trees, our experiments show where previous studies have claimed to produce optimal models yet failed; we show specific cases where this happens.

We demonstrate our methods on recidivism and credit risk datasets. These are two of the high-stakes decisions problems where interpretability is needed most in AI systems.

We provide ablation experiments to show which of our techniques is most influential over reducing computation for various datasets. As a result of this analysis, we are able to pinpoint possible future paths to improvement for scalability and computational speed. We believe these could be valuable for an emergent future class of possible new decision tree algorithms.

Author names are in alphabetic order.

ACM ISBN 000-0-0000-0000-0/00/00...\$15.00
<https://doi.org/00.0000/0000000.0000000>

2 RELATED WORK

Optimal decision trees have a quite a long history [3], so we focus on recent techniques and their relevance.

There are efficient algorithms that claim to generate optimal sparse trees, but which do not optimally balance the criteria of optimality and sparsity; instead they pre-specify the topology of the tree (*i.e.*, they know *a priori* exactly what the structure of the splits and leaves are, even though they do not know which variables are split) and only find the optimal tree of the given topology [16]. While this is an important problem, it is not the one we aim to solve, as we do not claim to know the topology of the optimal tree in advance. The most successful current algorithm of this variety is BinOCT [20], which searches for a full and complete binary tree of a given depth, based on continuous-optimization methods for continuous variables, which will be discussed at length later. Some exploration of learning optimal decision trees is based on boolean satisfiability (SAT) [17], but again, this work looks only for the optimal tree of a given number of nodes.

The DL8 algorithm [18] optimizes a ranking function to find a decision tree under constraints of size, depth, accuracy and leaves. DL8 creates trees from the bottom up, meaning that trees are assembled out of all possible leaves, which are itemsets pre-mined from the data [similarly to 2]. DL8 does not have publicly available sourcecode. Nijssen and Fromont [18] warn of issues of running out of memory when storing all partially-constructed trees.

There are several works that produce oblique trees [6], where splits involve several variables; oblique trees are not addressed here, as they can be less interpretable.

The most recent mathematical programming algorithms are OCT [5] and BinOCT [20]. Example figures from the OCT paper [5] show decision trees that are clearly suboptimal. However, as the code was not made public, the work in the OCT paper [5] is not easily reproducible, so it is not clear where the problem occurred. We discuss this in Section §5. Verwer and Zhang [20] introduce a much faster mathematical programming formulation, and their experiments indicate that BinOCT outperforms OCT, but since BinOCT is constrained to create complete binary trees of a given depth rather than optimally sparse trees, it sometimes creates unnecessary leaves in order to complete a tree at a given depth, as we show in Section §5. BinOCT solves a dramatically easier problem than the method introduced in this work. As it turns out, the search space of perfect binary trees of a given depth is much smaller than that of binary trees with the same number of leaves. For instance, the number of different unlabeled binary trees with 8 leaves is $Catalan(7) = 429$, but the number of unlabeled perfect binary trees with 8 leaves is only 1. In our setting, we penalize (but do not fix) the number of leaves, which means that our search space contains all trees, though we can bound the maximum number of leaves based on the size of the regularization parameter. Therefore, our search space is much larger than that of BinOCT.

Our work builds upon the CORELS algorithm of Angelino et al. [1, 2], Larus-Stone et al. [13] [and its predecessors 14, 21] which creates optimal decision lists (rule lists). Applying those ideas to decision trees is nontrivial. The rule list optimization problem is much easier, since the rules are pre-mined (there is no mining of rules in our decision tree optimization). Rule list optimization involves

Search Space of CORELS and Decision Trees

d	$p = 10$		$p = 20$	
	Rule Lists	Trees	Rule Lists	Trees
1	5.500×10^1	1.000×10^1	2.100×10^2	2.000×10^1
2	3.025×10^3	1.000×10^3	4.410×10^4	8.000×10^3
3	1.604×10^5	5.329×10^6	9.173×10^6	9.411×10^8
4	8.345×10^6	9.338×10^{20}	1.898×10^9	9.204×10^{28}
5	4.257×10^8	Inf	3.911×10^{11}	Inf

Table 1: Search spaces of rule lists and decision trees with number of variables $p = 10, 20$ and depth $d = 1, 2, 3, 4, 5$. The search space of the trees explodes in comparison.

selecting an optimal subset and an optimal permutation of the rules in each subset. Decision tree optimization involves considering every possible split of every possible variable, and every possible shape and size of tree. This is an exponentially harder optimization problem with a huge number of symmetries to consider. In addition, in CORELS, the maximum number of clauses per rule is set to be $c = 2$. For a data set with p features, there would be $D = p + \binom{p}{2}$ rules in total, and the number of distinct rule lists with d_r rules is $P(D, d_r)$, where $P(m, k)$ is the number of k -permutations of m . Therefore, the search space of CORELS is $\sum_{d_r=1}^{D-1} P(D, d_r)$. But, for a full binary tree with depth d_t and data with p features, the number of distinct trees is:

$$N_{d_t} = \sum_{n_0=1}^1 \sum_{n_1=1}^{2^{n_0}} \cdots \sum_{n_{d_t-1}=1}^{2^{n_{d_t-2}}} p \times \binom{2^{n_0}}{n_1} (p-1)^{n_1} \times \cdots \times \binom{2^{n_{d_t-2}}}{n_{d_t-1}} (p - (d_t - 1))^{n_{d_t-1}}, \quad (1)$$

and the search space of decision trees up to depth d is $\sum_{d_t=1}^d N_{d_t}$. Table 1 shows how the search spaces of rule lists and decision trees grow as the tree depth increases. The search space of the trees is massive compared to that of the rule lists.

Applying techniques from rule lists to decision trees necessitated new designs for the data structures, splitting mechanisms and bounds. An important difference between rule lists and trees is that during the growth of rule lists, we only add one new rule to the list each time, but for the growth of trees, we need to split existing leaves and add a new *pair* of leaves for each. This leads to several bounds that are quite different from those in CORELS, *i.e.*, Theorem 3.7, Theorem 3.8 and Corollary 3.9, which consider a pair of leaves rather than a single leaf. In this paper, we introduce bounds only for the case of one split at a time; however, in our implementation, we can split more than one leaf at a time, and the bounds are adapted accordingly.

3 OPTIMAL SPARSE DECISION TREES

We focus on binary classification, and our decision trees are Boolean functions. It is possible to generalize this framework to multiclass settings or weighted data. We denote training data as $\{(x_n, y_n)\}_{n=1}^N$, where $x_n \in \{0, 1\}^M$ are binary features and $y_n \in \{0, 1\}$ are labels. Let $\mathbf{x} = \{x_n\}_{n=1}^N$ and $\mathbf{y} = \{y_n\}_{n=1}^N$, and let $x_{n,m}$ denote the m -th feature of x_n .

For a decision tree, what we actually obtain is a set of leaves or conjunction of predicates. These leaves are mutually exclusive,

which means their order does not matter in evaluating the accuracy of the tree. A tree can thus be expressed in terms of only its leaves. A leaf set $d = (p_1, p_2, \dots, p_H)$ of length $H \geq 0$ is an H -tuple containing H distinct leaves, and every p_k corresponds to its $\hat{y}_k^{(\text{leaf})}$, for $k = 1, \dots, H$, where p_k is the classification rule of the path from the root to leaf k and $\hat{y}_k^{(\text{leaf})}$ is the label for leaf k . In our setting, p_k is a Boolean assertion, which evaluates to be either true or false for each datum x_n , and $\hat{y}_k^{(\text{leaf})}$ is a label prediction for all data in the leaf k . For example, $(x_{n,1} = 0) \wedge (x_{n,2} = 1)$ defines p_k and $(\hat{y}_n = 1)$ is the predicted label of datum x_n .

We explore the search space by considering which leaves of the tree can be beneficially split. The leaf set $d = (p_1, p_2, \dots, p_K, p_{K+1}, \dots, p_H)$ is the H -leaf tree, where the first K leaves are restricted not to be split, and the remaining $H - K$ leaves can be split. We alternately represent this leaf set as $d = (d_{un}, \delta_{un}, d_{\text{split}}, \delta_{\text{split}}, K, H)$, where $d_{un} = (p_1, \dots, p_K)$ are the unchanged leaves of d , $\delta_{un} = (\hat{y}_1^{(\text{leaf})}, \dots, \hat{y}_K^{(\text{leaf})}) \in \{0, 1\}^K$ are the predicted labels of leaves d_{un} , and $d_{\text{split}} = (p_{K+1}, \dots, p_H)$ are the leaves we are going to split, $\delta_{\text{split}} = (\hat{y}_{K+1}^{(\text{leaf})}, \dots, \hat{y}_H^{(\text{leaf})}) \in \{0, 1\}^{H-K}$ are the predicted labels of leaves d_{split} . We call d_{un} a K -prefix of d , which means its leaves are a size- K subset of $(p_1, \dots, p_K, \dots, p_H)$. If we have a new prefix d'_{un} , which is a supset of d_{un} , i.e., $d'_{un} \supseteq d_{un}$, then we say d'_{un} starts with d_{un} . We define $\sigma(d_{un})$ to be all d 's child trees whose (unchanged) leaves include d_{un} :

$$\sigma(d_{un}) = \{(d'_{un}, \delta'_{un}, d'_{\text{split}}, \delta'_{\text{split}}, K', H') : d'_{un} \text{ starts with } d_{un}\}. \quad (2)$$

If we have two prefix $d_{un} = (p_1, \dots, p_K)$ and $d'_{un} = (p_1, \dots, p_K, p_{K+1})$, and d'_{un} starts with d_{un} and contain one more leaf, then we define d'_{un} to be a child of d_{un} and d_{un} to be a parent of d'_{un} .

Note that in our setting, two trees with identical leaf sets, but with different assignments to d_{un} and d_{split} , are *different* trees. Further, a child tree can *only* be generated through splitting those leaves of its parent tree within d_{split} .

A tree d classifies datum x_n by providing the label prediction $\hat{y}_k^{(\text{leaf})}$ of the leaf whose p_k is true for x_n . If p_k evaluates to true for x_n , we say the leaf k of leaf set d_{un} captures x_n . In our setting, all the data captured by a prefix's leaves are also captured by the prefix itself.

Let β be a set of leaves. We define $\text{cap}(x_n, \beta) = 1$ if a leaf in β captures datum x_n , and 0 otherwise. For example, let d_{un} and d'_{un} be prefixes such that d'_{un} starts with d_{un} , then d'_{un} captures all the data that d_{un} captures: $\{x_n : \text{cap}(x_n, d_{un})\} \subseteq \{x_n : \text{cap}(x_n, d'_{un})\}$.

Now let d_{un} be a set of leaves, and let β be a subset of leaves in d_{un} . We denote the normalized support of β as $\text{supp}(\beta, x)$, which is the fraction of the dataset captured by β :

$$\text{supp}(\beta, x) = \frac{1}{N} \sum_{n=1}^N \text{cap}(x_n, \beta). \quad (3)$$

In our setting, decision trees are empirically set to minimize the regularized misclassification error. Therefore, the combination of training data (x, y) , a leaf set $d_{un} = (p_1, \dots, p_K)$, and $d_{\text{split}} = (p_{K+1}, \dots, p_H)$ denotes a tree $d = (d_{un}, \delta_{un}, d_{\text{split}}, \delta_{\text{split}}, K, H)$ with prefix d_{un} , where $\hat{y}_h^{(\text{leaf})}$, $(1 \leq h \leq H)$, corresponds to the majority label of data captured by leaf p_h within the set of label predictions $\delta_{un} =$

$(\hat{y}_1^{(\text{leaf})}, \dots, \hat{y}_K^{(\text{leaf})})$ and $\delta_{\text{split}} = (\hat{y}_{K+1}^{(\text{leaf})}, \dots, \hat{y}_H^{(\text{leaf})})$. Throughout this paper, we always assume a decision tree assigns labels to its leaves in this way.

3.1 Objective Function

For a tree $d = (d_{un}, \delta_{un}, d_{\text{split}}, \delta_{\text{split}}, K, H)$, we define its objective function as a combination of the misclassification error and a sparsity penalty on the number of leaves:

$$R(d, x, y) = \ell(d, x, y) + \lambda H(d). \quad (4)$$

$H(d)$ is the number of leaves in the tree d . $R(d, x, y)$ is a regularized empirical risk. The loss $\ell(d, x, y)$, is the misclassification error of d , i.e., the fraction of training data with incorrectly predicted labels. $\lambda H(d)$ is a regularization term which penalizes bigger trees. The regularization parameter $\lambda \geq 0$ is a small user-defined constant. It imposes a penalty whereby misclassifying a fraction of λ of the training data suffers the same as adding one to the number of leaves.

3.2 Optimization Framework

We minimize the objective function based on a branch-and-bound framework. We propose a series of specialized useful bounds that work together to eliminate a large part of the search space. These bounds are discussed in detail in the following paragraphs. Proofs are in the supplementary materials.

Some of our bounds could be adapted directly from CORELS [2], namely these two:

- **(Hierarchical objective lower bound)** Lower bounds of a parent tree also hold for every child tree of that parent. (§3.3, Theorem 3.1).
- **(Equivalent points bound)** For a given dataset, if there are multiple samples with exactly the same features but different labels, then no matter how we build our classifier, we would always make mistakes. And the lower bound of number of mistakes is the number of such samples with minority labels. (§3.7, Theorem 3.12).

Some of our bounds adapt from CORELS [1] with minor changes:

- **(Objective lower bound with one-step lookahead)** With respect to the number of leaves, if a tree copy does not achieve enough accuracy, we can prune all child trees of it. (§3.3, Lemma 3.2).
- **(A priori bound on the number of leaves)** For an optimal decision tree, we provide *a priori* an upper bound on the maximum number of leaves. (§3.4, Theorem 3.5).
- **(Lower bound on leaf support)** For an optimal tree, the support of each pair of leaves in the tree must be above 2λ . (§3.5, Theorem 3.6).

Some of our bounds are distinct from CORELS because they are only relevant to trees and not to lists:

- **(Lower bound on incremental classification accuracy)** Each split must result in sufficient reduction of the loss. The loss reduction should be larger than the increase of the regularization term. Otherwise, we need to further split at least one of the child leaves (§3.5, Theorem 3.7).
- **(Leaf permutation bound)** We only need to consider the optimal permutation of leaves in a tree; we do not need to consider all other permutations (explained in §3.6, Corollary 3.9).
- **(Leaf accurate support bound)** Each leaf in an optimal decision tree must have misclassification error that is sufficiently small. Specifically, for each leaf in an optimal decision tree, the number of correctly classified samples must be above a threshold. (§3.5, Theorem 3.8).

The supplement contains an additional set of bounds on the number of remaining tree evaluations.

3.3 Hierarchical Objective Lower Bound

The loss can be decomposed into two parts corresponding to the unchanged leaves and the leaves to be split:

$$\ell(d, \mathbf{x}, \mathbf{y}) \equiv \ell_p(d_{un}, \delta_{un}, \mathbf{x}, \mathbf{y}) + \ell_q(d_{split}, \delta_{split}, \mathbf{x}, \mathbf{y}),$$

where $d_{un} = (p_1, \dots, p_K)$, $\delta_{un} = (\hat{y}_1^{(leaf)}, \dots, \hat{y}_K^{(leaf)})$, $d_{split} = (p_{K+1}, \dots, p_H)$ and $\delta_{split} = (\hat{y}_{K+1}^{(leaf)}, \dots, \hat{y}_H^{(leaf)})$;

$$\ell_p(d_{un}, \delta_{un}, \mathbf{x}, \mathbf{y}) = \frac{1}{N} \sum_{n=1}^N \sum_{k=1}^K \text{cap}(x_n, p_k) \wedge \mathbb{1}[\hat{y}_k^{(leaf)} \neq y_n]$$

is the proportion of data in the unchanged leaves that are misclassified, and

$$\ell_q(d_{split}, \delta_{split}, \mathbf{x}, \mathbf{y}) = \frac{1}{N} \sum_{n=1}^N \sum_{k=K+1}^H \text{cap}(x_n, p_k) \wedge \mathbb{1}[\hat{y}_k^{(leaf)} \neq y_n]$$

is the proportion of data in the leaves we are going to split that are misclassified. We define a lower bound $b(d_{un}, \mathbf{x}, \mathbf{y})$ on the objective by leaving out the latter loss,

$$b(d_{un}, \mathbf{x}, \mathbf{y}) \equiv \ell_p(d_{un}, \delta_{un}, \mathbf{x}, \mathbf{y}) + \lambda H \leq R(d, \mathbf{x}, \mathbf{y}), \quad (5)$$

where the leaves d_{un} are kept and the leaves d_{split} are going to be split. Here, $b(d_{un}, \mathbf{x}, \mathbf{y})$ gives a lower bound on the objective of *any* child tree of d .

THEOREM 3.1 (HIERARCHICAL OBJECTIVE LOWER BOUND). *Define $b(d_{un}, \mathbf{x}, \mathbf{y}) = \ell_p(d_{un}, \delta_{un}, \mathbf{x}, \mathbf{y}) + \lambda H$, as in (5). Define $\sigma(d_{un})$ to be the set of all d 's child trees whose unchanged leaves contains d_{un} , as in (2). Let $d = (d_{un}, \delta_{un}, d_{split}, \delta_{split}, K, H)$ be a tree with unchanged leaves d_{un} , and let $d' = (d'_{un}, \delta'_{un}, d'_{split}, \delta'_{split}, K', H') \in \sigma(d_{un})$ be any child tree such that its unchanged leaves d'_{un} contain d_{un} and $K' \geq K, H' \geq H$, then $b(d_{un}, \mathbf{x}, \mathbf{y}) \leq R(d', \mathbf{x}, \mathbf{y})$.*

Consider a sequence of trees, where each tree is the parent of the following tree. In this case, the lower bounds of these trees increase monotonically, which is amenable to branch-and-bound. We illustrate our framework in Algorithm 1.

According to Theorem 3.1, we can hierarchically prune the search space. During the execution of the algorithm, we cache the current best (smallest) objective R^c , which is dynamic and monotonically decreasing. In this process, when we generate a tree whose unchanged leaves d_{un} corresponding to lower bound $b(d_{un}, \mathbf{x}, \mathbf{y}) \geq R^c$, according to Theorem 3.1, we do not need to consider *any* child tree $d' \in \sigma(d_{un})$ of this tree whose d'_{un} contains d_{un} . This is because we have $b(d'_{un}, \mathbf{x}, \mathbf{y}) \geq b(d_{un}, \mathbf{x}, \mathbf{y})$, which leads to $R(d', \mathbf{x}, \mathbf{y}) \geq b(d'_{un}, \mathbf{x}, \mathbf{y}) \geq b(d_{un}, \mathbf{x}, \mathbf{y}) \geq R^c$. Therefore, any child of such a tree cannot be the optimal tree, and we can prune the corresponding search space.

Based on Theorem 3.1, we describe a consequence in Lemma 3.2.

LEMMA 3.2 (OBJECTIVE LOWER BOUND WITH ONE-STEP LOOK-AHEAD). *Let d be a H -leaf tree with a K -leaf prefix and let R^c be the current best objective. If $b(d_{un}, \mathbf{x}, \mathbf{y}) + \lambda \geq R^c$, then for any child tree $d' \in \sigma(d_{un})$, its prefix d'_{un} starts with d_{un} and $K' > K, H' > H$, and it follows that $R(d', \mathbf{x}, \mathbf{y}) \geq R^c$.*

Algorithm 1 Branch-and-bound for learning optimal decision trees.

Input: Objective function $R(d, \mathbf{x}, \mathbf{y})$, objective lower bound $b(d_{un}, \mathbf{x}, \mathbf{y})$, set of features $S = \{s_m\}_{m=1}^M$, training data $(\mathbf{x}, \mathbf{y}) = \{(x_n, y_n)\}_{n=1}^N$, initial best known tree d^0 with objective $R^0 = R(d^0, \mathbf{x}, \mathbf{y})$; d^0 could be obtained as output from another (approximate) algorithm, otherwise, $(d^0, R^0) = (\text{null}, 1)$ provides reasonable default values

Output: Provably optimal decision tree d^* with minimum objective R^*

```

( $d^c, R^c$ )  $\leftarrow$  ( $d^0, R^0$ )            $\triangleright$  Initialize best tree and objective
 $Q \leftarrow \text{queue}([(((), (), \delta_{split}, 0, 0))])$   $\triangleright$  Initialize queue with empty tree
while  $Q$  not empty do            $\triangleright$  Stop when queue is empty
     $d = (d_{un}, \delta_{un}, d_{split}, \delta_{split}, K, H) \leftarrow Q.\text{pop}()$   $\triangleright$  Remove tree  $d$  from the queue
    if  $b(d_{un}, \mathbf{x}, \mathbf{y}) < R^c$  then    $\triangleright$  Bound: Apply Theorem 3.1
         $R \leftarrow R(d, \mathbf{x}, \mathbf{y})$         $\triangleright$  Compute objective of tree  $d$ 
        if  $R < R^c$  then          $\triangleright$  Update best tree and objective
            ( $d^c, R^c$ )  $\leftarrow$  ( $d, R$ )
        end if
        for every possible combination of features to split  $d_{split}$ 
             $\triangleright$  Branch: Enqueue  $d_{un}$ 's children
            split  $d_{split}$  and get new leaves  $d_{new}$ 
            for each possible subset  $d'_{split}$  of  $d_{new}$  do
                 $d'_{un} = d_{un} \cup (d_{new} \setminus d'_{split})$ 
                 $Q.\text{push}((d'_{un}, \delta'_{un}, d'_{split}, \delta'_{split}, K', H'))$ 
            end for
        end for
    end if
end while
( $d^*, R^*$ )  $\leftarrow$  ( $d^c, R^c$ )        $\triangleright$  Identify provably optimal solution

```

Therefore, even though we might have a tree with unchanged leaves d_{un} whose lower bound $b(d_{un}, \mathbf{x}, \mathbf{y}) \leq R^c$, if $b(d_{un}, \mathbf{x}, \mathbf{y}) + \lambda \geq R^c$, we can still prune all its child trees with a larger number of leaves and unchanged leaves $d'_{un} \supseteq d_{un}$.

3.4 Upper Bounds on Number of Leaves

During the branch-and-bound execution, the current best objective R^c implies an upper bound on the maximum number of leaves for those trees we still need to consider.

THEOREM 3.3 (UPPER BOUND ON THE NUMBER OF LEAVES). *For a dataset with M features, consider a state space of all trees. Let $L(d)$ be the number of leaves of tree d and let R^c be the current best objective. For all optimal trees $d^* \in \text{argmin}_d R(d, \mathbf{x}, \mathbf{y})$*

$$L(d^*) \leq \min \left(\left\lceil \frac{R^c}{\lambda} \right\rceil, 2^M \right), \quad (6)$$

where λ is the regularization parameter.

COROLLARY 3.4 (A PRIORI UPPER BOUND ON THE NUMBER OF LEAVES). *For all optimal tree $d^* \in \text{argmin}_d R(d, \mathbf{x}, \mathbf{y})$,*

$$L(d^*) \leq \min \left(\left\lceil \frac{1}{2\lambda} \right\rceil, 2^M \right). \quad (7)$$

For any particular tree d with unchanged leaves d_{un} , we can obtain potentially tighter upper bounds on the number of leaves for all its child trees whose unchanged leaves include d_{un} .

THEOREM 3.5 (PARENT-SPECIFIC UPPER BOUND ON THE NUMBER OF LEAVES). *Let $d = (d_{un}, \delta_{un}, d_{split}, \delta_{split}, K, H)$ be a tree, let $d' = (d'_{un}, \delta'_{un}, d'_{split}, \delta'_{split}, K', H')$ $\in \sigma(d_{un})$ be any child tree such that $d'_{un} \supseteq d_{un}$, and let R^c be the current best objective. If d'_{un} has lower bound $b(d'_{un}, \mathbf{x}, \mathbf{y}) < R^c$, then*

$$H' < \min \left(H + \left\lfloor \frac{R^c - b(d_{un}, \mathbf{x}, \mathbf{y})}{\lambda} \right\rfloor, 2^M \right). \quad (8)$$

Theorem 3.5 can be viewed as a generalization of the one-step lookahead bound (Lemma 3.2). This is because we can view (8) as a bound on $H' - H$, which provides an upper bound on the number of remaining splits we need to do based on the ideal tree we already have evaluated.

3.5 Lower Bounds on Leaf Support and Classification Accuracy

We provide three lower bounds on the fraction of correctly classified data and the normalized support of leaves in any optimal tree. All of them are relevant with λ .

THEOREM 3.6 (LOWER BOUND ON LEAF SUPPORT). *Let $d^* = (d_{un}, \delta_{un}, d_{split}, \delta_{split}, K, H)$ be any optimal tree with objective R^* , i.e., $d^* \in \arg\min_d R(d, \mathbf{x}, \mathbf{y})$. For each child leaf pair p_k, p_{k+1} of a split, the sum of normalized supports of p_k, p_{k+1} should be no less than twice the regularization parameter, i.e., 2λ ,*

$$2\lambda \leq \text{supp}(p_k, \mathbf{x}) + \text{supp}(p_{k+1}, \mathbf{x}). \quad (9)$$

Therefore, for a tree d , if any of its leaf pairs captures less than a fraction 2λ of samples, it is unlikely to be the optimal tree, even if $b(d_{un}, \mathbf{x}, \mathbf{y}) < R^*$. None of its child trees would be an optimal tree. Thus, after evaluating d , we can prune tree d . Theorem 3.6 only depends on the amount of data captured by the leaves. In Theorem 3.7, we provide a tighter upper bound on $R(d, \mathbf{x}, \mathbf{y})$ in (20) by also leveraging the classification accuracy.

THEOREM 3.7 (LOWER BOUND ON INCREMENTAL CLASSIFICATION ACCURACY). *Let $d^* = (d_{un}, \delta_{un}, d_{split}, \delta_{split}, K, H)$ be any optimal tree with objective R^* , i.e., $d^* \in \arg\min_d R(d, \mathbf{x}, \mathbf{y})$. Let d^* have leaves $d_{un} = (p_1, \dots, p_H)$ and labels $\delta_{un} = (\hat{y}_1^{(leaf)}, \dots, \hat{y}_H^{(leaf)})$. For each leaf pair p_k, p_{k+1} and the corresponding labels $\hat{y}_k^{(leaf)}, \hat{y}_{k+1}^{(leaf)}$ in d^* and the their parent node (the leaf in the parent tree) p_j and its label $\hat{y}_j^{(leaf)}$, define a_k to be the incremental classification accuracy of splitting p_j to get p_k, p_{k+1} :*

$$\begin{aligned} a_k \equiv & \frac{1}{N} \sum_{n=1}^N \{ \text{cap}(x_n, p_k) \wedge \mathbb{1}[\hat{y}_k^{(leaf)} = y_n] \\ & + \text{cap}(x_n, p_{k+1}) \wedge \mathbb{1}[\hat{y}_{k+1}^{(leaf)} = y_n] \\ & - \text{cap}(x_n, p_j) \wedge \mathbb{1}[\hat{y}_j^{(leaf)} = y_n] \}. \end{aligned} \quad (10)$$

In this case, λ provides a lower bound, $\lambda \leq a_k$.

Thus, when we split a leaf of the parent tree, if the incremental fraction of data that are correctly classified after this split is

less than a fraction λ , we need to further split at least one of the two child leaves to search for the optimal tree. We can actually view Theorem 3.6 as a sub-condition of Theorem 3.7, and in our implementation, we are able to leverage both of them. This is because Theorem 3.6 provides a much cheaper check in terms of the computation, which we can make use of to prune the search space before we calculate the bound of Theorem 3.7. In our algorithm, we apply Theorem 3.6 when we split the leaves. We only need to split the leaves whose normalized supports are no less than 2λ . And we apply Theorem 3.7 when we construct the trees. For every new split we do, we need to check the incremental accuracy corresponding to this split. If it is less than λ , we further split at least one of the two child leaves.

Both Theorem 3.6 and Theorem 3.7 are bounds for pair of leaves. We further give a bound on a single leaf's classification accuracy in Theorem 3.8.

THEOREM 3.8 (LOWER BOUND ON CLASSIFICATION ACCURACY). *Let $d^* = (d_{un}, \delta_{un}, d_{split}, \delta_{split}, K, H)$ be any optimal tree with objective R^* , i.e., $d^* \in \arg\min_d R(d, \mathbf{x}, \mathbf{y})$. For each leaf $(p_k, \hat{y}_k^{(leaf)})$ in d^* , the fraction of correctly classified data in leaf k should be no less than λ ,*

$$\lambda \leq \frac{1}{N} \sum_{n=1}^N \text{cap}(x_n, p_k) \wedge \mathbb{1}[\hat{y}_k^{(leaf)} = y_n]. \quad (11)$$

Thus, in a leaf we are thinking about extending, if a feature does not have the minimum fraction of correctly classified data for points that go to that leaf, then we can exclude that feature further down the tree extending that leaf. The feature will be excluded anywhere below the leaf, not just in the split immediately being considered on the current leaf.

3.6 Permutation Bound

If two trees are composed of the same leaves, i.e., they contain the same conjunctions of features up to a permutation, then they classify all the data in the same way and their child trees are also permutations of each other. Therefore, if we already have all children from one permutation of a tree, then there is no benefit to considering child trees generated from a different permutation.

COROLLARY 3.9 (LEAF PERMUTATION BOUND). *Let π be any permutation of $\{1, \dots, H\}$, Let $d = (d_{un}, \delta_{un}, d_{split}, \delta_{split}, K, H)$ and $D = (D_{un}, \Delta_{un}, D_{split}, \Delta_{split}, K, H)$ denote trees with leaves (p_1, \dots, p_H) and $D_{un} = (p_{\pi(1)}, \dots, p_{\pi(H)})$, respectively, i.e., the leaves in D correspond to a permutation of the leaves in d . Then the objective lower bounds of d and D are the same and their child trees correspond to permutations of each other.*

Therefore, if two trees have the same leaves, up to a permutation, according to Corollary 3.9, either of them can be pruned. We call this symmetry-aware pruning. In the next section, Section §3.6.1, we demonstrate how this help save many computations.

3.6.1 Upper bound on tree evaluations with symmetry-aware pruning. Here we give an upper bound on the total number of tree evaluations based on symmetry-aware pruning (§3.6). For every subset of K leaves, there are $K!$ leaf sets equivalent up to permutation. Thus, symmetry-aware pruning dramatically reduces the

search space by considering only one of them. This effects the execution of Algorithm 1's breadth-first search. With symmetry-aware pruning, when it evaluates trees of size K , for each set of trees equivalent up to a permutation, it keeps only a single tree.

THEOREM 3.10 (UPPER BOUND ON TREE EVALUATIONS WITH SYMMETRY-AWARE PRUNING). *Consider a state space of all trees formed from a set S of 3^M leaves where M is the number of features (the 3 options correspond to having the feature's value be 1, having its value be 0, or not including the feature), and consider the branch-and-bound algorithm with symmetry-aware pruning. Define $\Gamma_{\text{tot}}(S)$ to be the total number of prefixes evaluated. For any set S of 3^M leaves,*

$$\Gamma_{\text{tot}}(S) \leq 1 + \sum_{k=1}^K N_k + C(M, k) - P(M, k), \quad (12)$$

where $K = \min(\lfloor 1/2\lambda \rfloor, 2^M)$, N_k is defined in (1).

PROOF. By Corollary 3.4, $K \equiv \min(\lfloor 1/2\lambda \rfloor, 2^M)$ gives an upper bound on the number of leaves of any optimal tree. The algorithm begins by evaluating the empty tree, followed by M trees of depth $k = 1$, then $N_2 = \sum_{n_0=1}^1 \sum_{n_1=1}^{2^{n_0}} M \times \binom{2^{n_0}}{n_1} (M-1)^{n_1}$ trees of depth $k = 2$. Before proceeding to length $k = 3$, we keep only $N_2 + C(M, 2) - P(M, 2)$ trees of depth $k = 2$, where N_k is defined in (1), $P(M, k)$ denotes the number of k -permutations of M and $C(M, k)$ denotes the number of k -combinations of M . Now, the number of length $k = 3$ prefixes we evaluate is $N_3 + C(M, 3) - P(M, 3)$. Propagating this forward gives (12). \square

Pruning based on permutation symmetries thus yields significant computational savings of $\sum_{k=1}^K P(M, k) - C(M, k)$. For example, when $M = 10$ and $K = 5$, the number reduced due to symmetry-aware pruning is about 35463. If $M = 20$ and $K = 10$, the number of evaluations is reduced by about 7.36891×10^{11} .

3.7 Equivalent Points Bound

Equivalent Points Bound applies when multiple observations captured by a leaf in d_{split} have identical features but opposite labels: then no tree, including those that extend d_{split} , can correctly classify all of these observations. In this case, the number of misclassifications made by any tree is no less than the number of observations with the minority label.

For a data set $\{(x_n, y_n)\}_{n=1}^N$ and also a set of features $\{s_m\}_{m=1}^M$, we define a set of samples to be equivalent if they have exactly the same feature values, i.e., $x_i \neq x_j$ are equivalent if

$$\frac{1}{M} \sum_{m=1}^M \mathbb{1}[\text{cap}(x_i, s_m) = \text{cap}(x_j, s_m)] = 1.$$

Note that a data set consists of multiple sets of equivalent points; let $\{e_u\}_{u=1}^U$ enumerate these sets. For each observation x_i , it belongs to a equivalent points set e_u . We denote the fraction of data with the minority label in set e_u as $\theta(e_u)$, e.g., let $e_u = \{x_n : \forall m \in [M], \mathbb{1}[\text{cap}(x_n, s_m) = \text{cap}(x_i, s_m)]\}$, and let q_u be the minority class label among points in e_u , then

$$\theta(e_u) = \frac{1}{N} \sum_{n=1}^N \mathbb{1}[x_n \in e_u] \mathbb{1}[y_n = q_u]. \quad (13)$$

We can combine the equivalent points bound with other bounds to get a tighter lower bound on the objective function. As the experimental results demonstrate in §5, there is sometimes a substantial reduction of the search space after incorporating the equivalent points bound. We propose a general equivalent points bound in Proposition 3.11. We incorporate it into our framework by proposing the specific equivalent points bound in Theorem 3.12.

PROPOSITION 3.11 (GENERAL EQUIVALENT POINTS BOUND). *Let $d = (d_{\text{un}}, \delta_{\text{un}}, d_{\text{split}}, \delta_{\text{split}}, K, H)$ be a tree, then*

$$R(d, \mathbf{x}, \mathbf{y}) \geq \sum_{u=1}^U \theta(e_u) + \lambda H.$$

Recall that in our lower bound $b(d_{\text{un}}, \mathbf{x}, \mathbf{y})$ in (5), we leave out the misclassification errors of leaves we are going to split $\ell_0(d_{\text{split}}, \delta_{\text{split}}, \mathbf{x}, \mathbf{y})$ from the objective $R(d, \mathbf{x}, \mathbf{y})$. Incorporating the equivalent points bound in Theorem 3.12, we obtain a tighter bound on our objective because we now have a tighter lower bound on the misclassification errors of leaves we are going to split, $0 \leq b_0(d_{\text{split}}, \mathbf{x}, \mathbf{y}) \leq \ell_0(d_{\text{split}}, \delta_{\text{split}}, \mathbf{x}, \mathbf{y})$.

THEOREM 3.12 (EQUIVALENT POINTS BOUND). *Let d be a tree with leaves $d_{\text{un}}, d_{\text{split}}$ and lower bound $b(d_{\text{un}}, \mathbf{x}, \mathbf{y})$, then for any tree $d' \in \sigma(d)$ whose prefix $d'_{\text{un}} \supseteq d_{\text{un}}$,*

$$R(d', \mathbf{x}, \mathbf{y}) \geq b(d_{\text{un}}, \mathbf{x}, \mathbf{y}) + b_0(d_{\text{split}}, \mathbf{x}, \mathbf{y}), \quad \text{where} \quad (14)$$

$$b_0(d_{\text{split}}, \mathbf{x}, \mathbf{y}) = \frac{1}{N} \sum_{u=1}^U \sum_{n=1}^N \text{cap}(x_n, d_{\text{split}}) \wedge \mathbb{1}[x_n \in e_u] \mathbb{1}[y_n = q_u]. \quad (15)$$

4 IMPLEMENTATION

Much of our implementation effort revolves around exploiting incremental computation. During the execution of Algorithm 1, for each tree d , we compute the lower bound $b(d_{\text{un}}, \mathbf{x}, \mathbf{y})$ of the tree based on its unchanged leaves d_{un} and the corresponding objective $R(d, \mathbf{x}, \mathbf{y})$ of the tree. Given the hierarchical nature of the parent-children relationship, we *incrementally* compute the objective function and the lower bound throughout the branch-and-bound execution of the algorithm. Together, these ideas save >97% execution time.

We implement a series of data structures designed to support this incremental computation. First, we store bounds and intermediate results for the trees and the leaves to support incremental computation of the lower bound and the objective. As a global statistic, we maintain the best (minimum) observed value of the objective function and the corresponding tree. As each leaf in a tree represents a set of clauses, for each leaf we store, in an efficient bit-vector representation, the set of samples captured by that clause and the accuracy of those samples with respect to the prediction of the leaf. From the values in the leaves, we can efficiently compute both the value of the objective for an entire tree and leaf values for children formed from splitting the leaf.

Second, we use a priority queue to order the exploration of the search space. The queue serves as a worklist with each entry in the queue corresponding to a tree. When an entry is removed from the queue, we use it to generate child trees, incrementally computing the information for the child trees. The ordering of the

worklist represents a scheduling policy. We evaluated both structural orderings, e.g., breadth first search and depth first search, and metric-based orderings, e.g., objective function, the lower bound. Each metric produces a different scheduling policy. We achieve the best performance in runtime and memory consumption using the curiosity metric from CORELS [2], which is the objective lower bound, divided by the normalized support of its unchanged leaves. For example, relative to using the objective, curiosity reduces runtime by a factor of two and memory consumption by a factor of four.

Third, to support symmetry-aware pruning from Corollary 3.9, we introduce two symmetry-aware maps – a LeafCache and a TreeCache. A leaf is a set of clauses, each of which corresponds to an attribute and the value (0,1) of that attribute. As the leaves of a decision tree are mutually exclusive, the data captured by each leaf is insensitive to the order of the leaf’s clauses. We encode leaves in a canonical order (sorted by attribute indexes) and use that canonical order as the key into the LeafCache. Each entry in the LeafCache represents all permutations of a set of clauses. We use a Python dictionary to map these keys to the leaf and its cached values. Before creating a leaf object, we first check whether or not we have already created an identical one. If we have not, then we create the leaf and insert it into the map. Otherwise, the permutation already exists, so instead of creating it, we retrieve it from the map and use it as the leaf of the tree we are constructing. Next, we implement the permutation bound (Corollary 3.9) using the TreeCache. The TreeCache contains encodings of all the trees we have evaluated. Like we did for the clauses in the LeafCache, we introduce a canonical order over the leaves in a tree and use that as the key to the TreeCache. If our algorithm produces a tree that is simply a permutation of a tree we have already evaluated, then we need not evaluate it again. We determine that this is the case by looking up the tree under consideration in the TreeCache. If it’s in the cache, we do nothing; if it is not, then we compute the bounds for the tree and insert it into the cache.

Now, we illustrate how these data structures support execution of our algorithm. We initialize the algorithm with the current best objective R^c and tree d^c . For unexplored trees in the queue, the scheduling policy selects the next tree d to split; we keep removing elements from the queue until the queue is empty. Then, for every very possible combination of features to split d_{split} , we construct a new tree d' with incremental calculation of the lower bound $b(d'_{un}, \mathbf{x}, \mathbf{y})$ and the objective $R(d', \mathbf{x}, \mathbf{y})$. If we achieve a better objective $R(d', \mathbf{x}, \mathbf{y})$, i.e., less than the current best objective R^c , we update R^c and d^c . If the lower bound of the new tree d' , combined with Theorem 3.12 and Theorem 3.2, is less than the current best objective, then we push it into the queue. Otherwise, according to the hierarchical lower bound (Theorem 3.1), no child of d' could possibly have an objective better than R^c , which means we do not push d' queue. When there are no more trees to explore, i.e., the queue is empty, we have finished the search of the whole space and output the (provably) optimal tree.

5 EXPERIMENTS

We address the following questions through experimental analysis: (1) *Do existing methods achieve optimal solutions, and if not, how*

far are they from optimal? (2) *How fast does our algorithm converge given the hardness of the problem it is solving?* (3) *How much does each of the bounds contribute to the performance of our algorithm?* (4) *What do optimal trees look like?*

The results of per-bound performance and memory improvement (Table 2) were run on a *m5a.4xlarge* instance of AWS’s Elastic Compute Cloud (EC2). The instance has 16 2.5GHz virtual GPU cores and 64 GB of RAM. All other results were run on a personal laptop with a 2.4GHz i5-8259U processor and 16GB of RAM.

We used 7 datasets: Five of them are from the UCI Machine Learning Repository [8], (tic-tac-toe, car evaluation, monk1, monk2, monk3). The other two datasets are the ProPublica recidivism data set [12] and the Fair Isaac (FICO) credit risk datasets [9]. We predict which individuals are arrested within two years of release ($N = 7,215$) on the recidivism data set and whether an individual will default on a loan for the FICO dataset.

Accuracy and optimality: We tested the accuracy of our algorithm against baseline methods CART and BinOCT [20]. BinOCT is the most recent publically available method for learning optimal classification trees and was shown to outperform other previous methods. As far as we know, there is no public code for most of the other relevant baselines, including [5, 6, 16]. One of these methods, OCT [5], reports that CART often dominates their performance (see Fig. 4 and Fig. 5 in their paper). Our models can never be worse than CART’s models, because in our implementation, we use the objective value of CART’s solution as a warm start to the objective value of the current best.

Figure 1 shows the training accuracy on each dataset. The time limits for both BinOCT and our algorithm are set to be 30 minutes.

Main results: We observe that (i) We can now evaluate how close to optimal other methods are (and they are often close to optimal or optimal). (ii) Sometimes, the baselines are *not* optimal. Recall that BinOCT searches only for the optimal tree *given the topology of the complete binary tree of a certain depth*. This restriction on the topology massively reduces the search space so that BinOCT runs quickly, but in exchange, it misses optimal sparse solutions that our method finds. (iii) Our method is *fast*. Our method runs on only one thread (we have not yet parallelized it) whereas BinOCT is highly optimized; it makes use of all eight threads in our experiments. Even with the $\approx 8x$ speedup of BinOCT, our method is still competitive.

Convergence: Figure 2 illustrates the behavior of OSDT for the ProPublica COMPAS dataset with $\lambda = 0.005$, for two different scheduling policies (curiosity and lower bound). The charts in the figure show how the current best objective value R^c and the lower bound $b(d_{un}, \mathbf{x}, \mathbf{y})$ vary as the algorithm progresses. When we schedule using the lower bound, the lower bounds of evaluated trees increase monotonically, and OSDT certifies optimality only when the value of the lower bound becomes large enough that we can prune the remaining search space or when the queue is empty, whichever is reached earlier. Using the curiosity, OSDT finds the optimal tree much more quickly than when using the lower bound.

Scalability: Figure 3 shows the scalability of OSDT with respect to the number of samples and the number of features. Running time grows *linearly* in the number of samples (n), allowing us to scale to millions of samples. With a small number of features (i.e., 4-8) runtime grows exponentially, however, as we add more features, the growth rate slows down, because we are able to prune the search

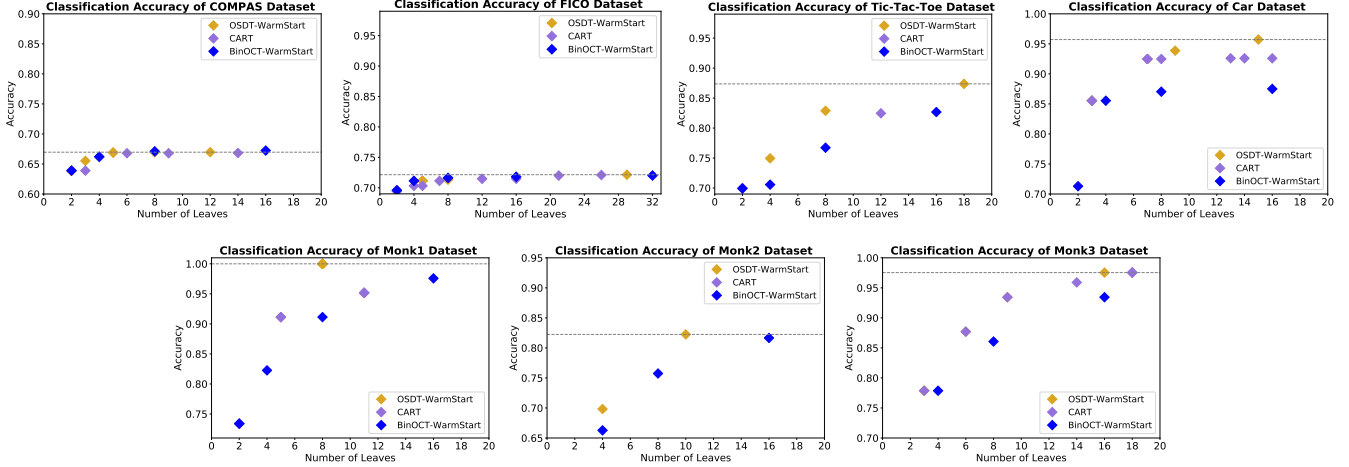


Figure 1: Training accuracy of OSDT, CART, BinOCT on different data (time limit: 30min). Horizontal lines indicate the accuracy of the best OSDT tree. On most datasets, all trees of BinOCT and CART are below this line.

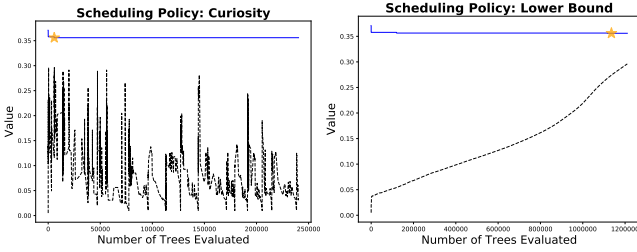


Figure 2: Example OSDT execution traces (COMPAS Dataset, $\lambda = 0.005$). Lines are the objective value and dashes are the lower bound for OSDT. For each scheduling policy, we mark the time to optimum and the optimal objective value using a star.

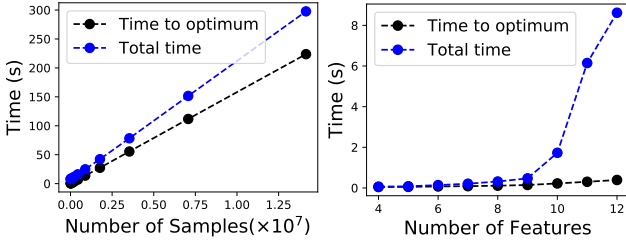


Figure 3: Scalability with respect to number of samples and number of features using (multiples of) the ProPublica data set. ($\lambda = 0.005$).

space more quickly as the number of features grows. For example, it takes us half the runtime to find the optimal tree using four features, but only one quarter of the run time to find the optimal tree with twelve features.

Algorithm optimizations: Next, we evaluate how much each of our bounds contributes to OSDT’s performance and what effect the scheduling metric has on execution. Table 2 provides experimental statistics of total execution time, time to optimum, total number of trees evaluated, number of trees evaluated to optimum, and memory consumption on the recidivism data set. The first row is the full OSDT implementation, and the others are variants each of which removes a specific bound. While all the optimizations reduce

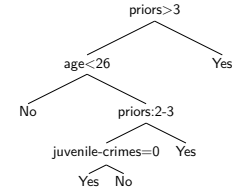


Figure 4: The optimal decision tree generated by OSDT on COMPAS dataset. ($\lambda = 0.005$)

the search space, the lookahead and equivalent points bounds are, by far, the most significant, reducing time to optimum by at least two orders of magnitude and reducing memory consumption by more than one order of magnitude. In our experiment, although the scheduling policy has a smaller effect, it is still significant – curiosity is a factor of two faster than the objective function and consumes 25% of the memory consumed when using the objective function. All other scheduling policies, *i.e.*, the lower bound and the entropy, are significantly worse.

Trees: Finally, we provide illustrations of the trees produced by OSDT and the baseline methods in Figures 4, 5 and 6. OSDT generates trees of any shape, and our objective penalizes trees with more leaves, thus it never introduces splits that produce a pair of leaves with the same label. However, the trees generated by BinOCT are always complete binary trees of given depth. Limiting the tree shape prevents BinOCT from finding the globally optimal tree in most cases. In fact, BinOCT sometimes produces useless splits, which leads to a trees with more leaves than necessary to achieve the same accuracy.

6 CONCLUSION AND FUTURE WORK

Our work shows the possibility of optimal (or provably near-optimal) decision trees for reasonably sized datasets. We have reason to believe this basic framework can be extended to much larger datasets with some more work. Angelino et al. [2] identified a key mechanism for scaling these algorithms up, which is the possibility of a bound stating that highly correlated features can substitute for each

Per-bound performance improvement (ProPublica data set)

Algorithm variant	Total time (s)	Slow-down	Time to optimum (s)
All bounds	14.75	—	1.09
No support bound	16.69	1.13×	1.18
No incremental accuracy bound	29.42	1.99×	1.27
No accuracy bound	31.72	2.15×	1.44
No lookahead bound	31479	2134×	186
No equivalent points bound	>8226	>557×	—

Algorithm variant	Total #trees evaluated	#trees to optimum	Mem (GB)
All bounds	241306	18195	.08
No support bound	278868	21232	.08
No incremental accuracy bound	548618	24682	.08
No accuracy bound	482013	23075	.09
No lookahead bound	283712499	3078445	10
No equivalent points bound	>41000000	—	>64

Table 2: Per-bound performance improvement, for the ProPublica data set ($\lambda = 0.005$, cold start, using curiosity). The columns report the total execution time, time to optimum, total number of trees evaluated, number of trees evaluated to optimum, and memory consumption. The first row shows our algorithm with all bounds; subsequent rows show variants that each remove a specific bound (one bound at a time, not cumulative). All rows except the last one represent a complete execution, i.e., until the queue is empty. For the last row ('No equivalent points bound'), the algorithm was terminated after running out of the memory (about ~64GB RAM).

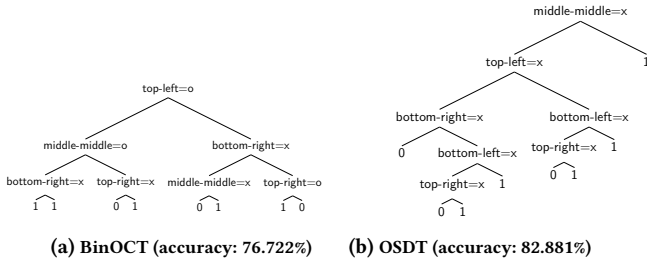


Figure 5: The decision tree generated by BinOCT and OSDT on the Tic-Tac-Toe data. Trees of BinOCT must be complete binary trees, while OSDT can generate trees of any shape.

other, leading to similar model accuracies. We have tried to put this bound into OSDT, but the amount of time spent evaluating the bound increased computation rather than reduce it. Further work needs to be done to determine when the bound should be computed in order to be effective, but if it can be done, the algorithm could potentially scale to much larger scales.

AVAILABILITY

The code is publicly available at <removed to preserve anonymity>.

REFERENCES

- [1] E. Angelino, N. Larus-Stone, D. Alabi, M. Seltzer, and C. Rudin. 2017. Learning certifiably optimal rule lists for categorical data. In *ACM SIGKDD International*

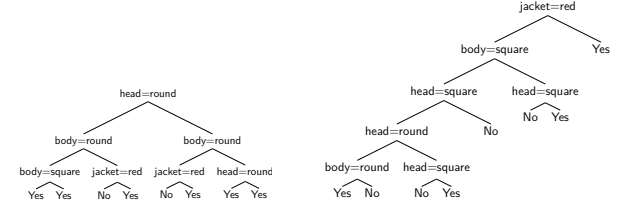


Figure 6: The decision tree generated by BinOCT and OSDT on Monk1 dataset. The tree generated by BinOCT includes useless splits, while OSDT can avoid this problem.

Conference on Knowledge Discovery and Data Mining (KDD).

- [2] Elaine Angelino, Nicholas Larus-Stone, Daniel Alabi, Margo Seltzer, and Cynthia Rudin. 2018. Learning Certifiably Optimal Rule Lists for Categorical Data. *Journal of Machine Learning Research* 18, 234 (2018), 1–78. <http://jmlr.org/papers/v18/17-716.html>
- [3] Kristin Bennett. 1992. Decision tree construction via linear programming. In *Proceedings of the 4th Midwest Artificial Intelligence and Cognitive Science Society Conference, Utica, Illinois*.
- [4] K. P. Bennett and J. A. Blue. 1996. *Optimal Decision Trees*. Technical Report. R.P.I. Math Report No. 214, Rensselaer Polytechnic Institute.
- [5] Dimitris Bertsimas and Jack Dunn. 2017. Optimal classification trees. *Machine Learning* 106, 7 (2017), 1039–1082.
- [6] Rafael Blanquero, Emilio Carrizosa, Cristina Molero-Rio, and Dolores Romero Morales. 2018. Optimal Randomized Classification Trees. (Aug. 2018).
- [7] L. Breiman, J. H. Friedman, R. A. Olshen, and C. J. Stone. 1984. *Classification and Regression Trees*. Wadsworth.
- [8] Dua Dheeru and Efi Karra Taniskidou. 2017. UCI Machine Learning Repository. <http://archive.ics.uci.edu/ml>
- [9] FICO, Google, Imperial College London, MIT, University of Oxford, UC Irvine, and UC Berkeley. 2018. FICO Explainable Machine Learning Challenge. <https://community.fico.com/s/explainable-machine-learning-challenge>
- [10] Anthony W. Flores, Christopher T. Lowenkamp, and Kristin Bechtel. 2016. False Positives, False Negatives, and False Analyses: A Rejoinder to “Machine Bias: There’s Software Used Across the Country to Predict Future Criminals”. *Federal probation* 80, 2 (September 2016).
- [11] Adam R. Klivans and Rocco A. Servedio. 2006. Toward Attribute Efficient Learning of Decision Lists and Parities. *Journal of Machine Learning Research* 7 (2006), 587–602.
- [12] J. Larson, S. Mattu, L. Kirchner, and J. Angwin. 2016. How We Analyzed the COMPAS Recidivism Algorithm. *ProPublica* (2016).
- [13] N. Larus-Stone, E. Angelino, D. Alabi, M. Seltzer, V. Kaxiras, A. Saligrama, and C. Rudin. 2018. Systems Optimizations for Learning Certifiably Optimal Rule Lists. In *SysML Conference*.
- [14] B. Letham, C. Rudin, T. H. McCormick, and D. Madigan. 2015. Interpretable classifiers using rules and Bayesian analysis: Building a better stroke prediction model. *The Annals of Applied Statistics* 9, 3 (2015), 1350–1371.
- [15] Michael McGough. 2018. How bad is Sacramento’s air, exactly? Google results appear at odds with reality, some say. *Sacramento Bee* (2018). <https://www.sacbee.com/news/state/california/fires/article216227775.html>
- [16] Matt Menickelly, Oktay Günlük, Jayant Kalagnanam, and Katya Scheinberg. 2018. Optimal Decision Trees for Categorical Data via Integer Programming. *Preprint at arXiv:1612.03225* (Jan. 2018). <http://arxiv.org/abs/1612.03225>
- [17] Nina Narodytska, Alexey Ignatiev, Filipe Pereira, Joao Marques-Silva, and IS-DCT SB RAS. 2018. Learning Optimal Decision Trees with SAT. In *IJCAL* 1362–1368.
- [18] Siegfried Nijssen and Elisa Fromont. 2007. Mining optimal decision trees from itemset lattices. In *Proceedings of the 13th ACM SIGKDD international conference on Knowledge discovery and data mining*. ACM, 530–539.
- [19] J. R. Quinlan. 1993. *C4.5: Programs for Machine Learning*. Morgan Kaufmann.
- [20] Sizzo Verwer and Yingqian Zhang. 2019. Learning optimal classification trees using a binary linear program formulation. In *33rd AAAI Conference on Artificial Intelligence*.
- [21] H. Yang, C. Rudin, and M. Seltzer. 2017. Scalable Bayesian Rule Lists. In *International Conference on Machine Learning (ICML)*.
- [22] John R. Zech, Marcus A Badgeley, Manway Liu, Anthony B. Costa, Joseph J. Titano, and Eric K. Oermann. 2018. Confounding Variables Can Degrade Generalization Performance of Radiological Deep Models. *arXiv:1807.00431* (2018).

A UPPER BOUNDS ON NUMBER OF TREE EVALS

In this section, based on the upper bounds on the number of leaves from §3.4, we give corresponding upper bounds on the number of tree evaluations made by Algorithm 1. First, in Theorem A.1, based on information about the state of Algorithm 1's execution, we calculate, for any given execution state, upper bounds on the number of additional tree evaluations needed for the execution to complete. We define the number of *remaining tree evaluations* as the number of trees that are currently in, or will be inserted into, the queue. We evaluate the number of tree evaluations based on current execution information of the current best objective R^c and the trees in the queue Q of Algorithm 1.

THEOREM A.1 (UPPER BOUND ON NUMBER OF REMAINING TREE EVALUATIONS). *Consider the state space of all possible leaves formed from a set of M features, and consider Algorithm 1 at a particular instant during execution. Denote the current best objective as R^c , the queue as Q , and the length of prefix d_{un} as $L(d_{un})$. Denoting the number of remaining prefix evaluations as $\Gamma(R^c, Q)$, the bound is:*

$$\Gamma(R^c, Q) \leq \sum_{d_{un} \in Q} \sum_{k=0}^{f(d_{un})} \frac{(3^M - L(d_{un}))!}{(3^M - L(d_{un}) - k)!}, \quad (16)$$

$$\text{where } f(d_{un}) = \min \left(\left\lfloor \frac{R^c - b(d_{un}, \mathbf{x}, \mathbf{y})}{\lambda} \right\rfloor, 3^M - L(d_{un}) \right). \quad (17)$$

The corollary below is a naïve upper bound on the total number of tree evaluations during the process of Algorithm 1's execution. It does not use algorithm execution state to bound the size of the search space like Theorem A.1, and it relies only on the number of features and the regularization parameter λ .

COROLLARY A.2 (UPPER BOUND ON THE TOTAL NUMBER OF TREE EVALUATIONS). *Define $\Gamma_{\text{tot}}(S)$ to be the total number of tree evaluated by Algorithm 1, given the state space of all possible leaves formed from a set S of M features. For any set S of all leaves formed of M features,*

$$\Gamma_{\text{tot}}(S) \leq \sum_{k=0}^K \frac{3^M!}{(3^M - k)!}, \text{ where } K = \min(\lfloor 1/2\lambda \rfloor, 2^M).$$

B PROOF OF THEOREMS

B.1 Proof of Theorem 3.3

PROOF. For an optimal tree d^* with objective R^* ,

$$\lambda L(d^*) \leq R^* = R(d^*, \mathbf{x}, \mathbf{y}) = \ell(d^*, \mathbf{x}, \mathbf{y}) + \lambda L(d^*) \leq R^c.$$

The maximum possible number of leaves for d^* occurs when $\ell(d^*, \mathbf{x}, \mathbf{y})$ is minimized; therefore this gives bound (6).

For the rest of the proof, let $H^* = L(d^*)$ be the length of d^* . If the current best tree d^c has zero misclassification error, then

$$\lambda H^* \leq \ell(d^*, \mathbf{x}, \mathbf{y}) + \lambda H^* = R(d^*, \mathbf{x}, \mathbf{y}) \leq R^c = R(d^c, \mathbf{x}, \mathbf{y}) = \lambda H,$$

and thus $H^* \leq H$. If the current best tree is suboptimal, i.e., $d^c \notin \text{argmin}_d R(d, \mathbf{x}, \mathbf{y})$, then

$$\lambda H^* \leq \ell(d^*, \mathbf{x}, \mathbf{y}) + \lambda H^* = R(d^*, \mathbf{x}, \mathbf{y}) < R^c = R(d^c, \mathbf{x}, \mathbf{y}) = \lambda H,$$

in which case $H^* < H$, i.e., $H^* \leq H - 1$, since H is an integer. \square

B.2 Proof of Theorem 3.5

PROOF. First, note that $H' \geq H$. Now recall that

$$\begin{aligned} b(d_{un}, \mathbf{x}, \mathbf{y}) &= \ell_p(d_{un}, \delta_{un}, \mathbf{x}, \mathbf{y}) + \lambda H \\ &\leq \ell_p(d'_{un}, \delta'_{un}, \mathbf{x}, \mathbf{y}) + \lambda H' = b(d'_{un}, \mathbf{x}, \mathbf{y}), \end{aligned}$$

and that $\ell_p(d_{un}, \delta_{un}, \mathbf{x}, \mathbf{y}) \leq \ell_p(d'_{un}, \delta'_{un}, \mathbf{x}, \mathbf{y})$. Combining these bounds and rearranging gives

$$\begin{aligned} b(d'_{un}, \mathbf{x}, \mathbf{y}) &= \ell_p(d'_{un}, \delta'_{un}, \mathbf{x}, \mathbf{y}) + \lambda H + \lambda(H' - H) \\ &\geq \ell_p(d_{un}, \delta_{un}, \mathbf{x}, \mathbf{y}) + \lambda H + \lambda(H' - H) \\ &= b(d_{un}, \mathbf{x}, \mathbf{y}) + \lambda(H' - H). \end{aligned} \quad (18)$$

Combining (??) with $b(d'_{un}, \mathbf{x}, \mathbf{y}) < R^c$ gives (8). \square

B.3 Proof of Theorem A.1

PROOF. The number of remaining tree evaluations is equal to the number of trees that are currently in or will be inserted into queue Q . For any such tree with unchanged leaves d_{un} , Theorem 3.5 gives an upper bound on the length of tree with unchanged leaves d'_{un} that contains d_{un} :

$$L(d'_{un}) \leq \min \left(L(d_{un}) + \left\lfloor \frac{R^c - b(d_{un}, \mathbf{x}, \mathbf{y})}{\lambda} \right\rfloor, 2^M \right) \equiv U(d_{un}).$$

This gives an upper bound on the remaining tree evaluations:

$$\begin{aligned} \Gamma(R^c, Q) &\leq \sum_{d_{un} \in Q} \sum_{k=0}^{U(d_{un}) - L(d_{un})} P(3^M - L(d_{un}), k) \\ &= \sum_{d_{un} \in Q} \sum_{k=0}^{f(d_{un})} \frac{(3^M - L(d_{un}))!}{(3^M - L(d_{un}) - k)!}, \end{aligned} \quad (19)$$

where $P(m, k)$ denotes the number of k -permutations of m . \square

B.4 Proof of Proposition A.2

PROOF. By Corollary 3.4, $K \equiv \min(\lfloor 1/2\lambda \rfloor, 2^M)$ gives an upper bound on the number of leaves of any optimal tree. Since we can think of our problem as finding the optimal selection and permutation of k out of 3^M leaves, over all $k \leq K$,

$$\Gamma_{\text{tot}}(S) \leq 1 + \sum_{k=1}^K P(3^M, k) = \sum_{k=0}^K \frac{3^M!}{(3^M - k)!}.$$

\square

B.5 Proof of Theorem 3.6

PROOF. Let $d^* = (d_{un}, \delta_{un}, d_{\text{split}}, \delta_{\text{split}}, K, H)$ be an optimal tree with leaves (p_1, \dots, p_H) and labels $(\hat{y}_1^{(\text{leaf})}, \dots, \hat{y}_H^{(\text{leaf})})$. Consider the tree $d = (d'_{un}, \delta'_{un}, d'_{\text{split}}, \delta'_{\text{split}}, K', H')$ derived from d^* by deleting a pair of leaves $p_i \rightarrow \hat{y}_i^{(\text{leaf})}$, $p_{i+1} \rightarrow \hat{y}_{i+1}^{(\text{leaf})}$ and adding the their parent leaf $p_j \rightarrow \hat{y}_j^{(\text{leaf})}$, therefore $d'_{un} = (p_1, \dots, p_{i-1}, p_{i+2}, \dots, p_H, p_j)$ and $\delta'_{un} = (\hat{y}_1^{(\text{leaf})}, \dots, \hat{y}_{i-1}^{(\text{leaf})}, \hat{y}_{i+2}^{(\text{leaf})}, \dots, \hat{y}_H^{(\text{leaf})}, \hat{y}_j^{(\text{leaf})})$.

When d misclassified half of the data captured by p_i, p_{i+1} , while d^* correctly classified them all, the difference between d and d^*

would maximize, which provides an upper bound:

$$\begin{aligned}
R(d, \mathbf{x}, \mathbf{y}) &= \ell(d, \mathbf{x}, \mathbf{y}) + \lambda(H - 1) \\
&\leq \ell(d^*, \mathbf{x}, \mathbf{y}) + \text{supp}(p_i, \mathbf{x}) + \text{supp}(p_{i+1}, \mathbf{x}) \\
&\quad - \frac{1}{2}[\text{supp}(p_i, \mathbf{x}) + \text{supp}(p_{i+1}, \mathbf{x})] + \lambda(H - 1) \\
&= R(d^*, \mathbf{x}, \mathbf{y}) + \frac{1}{2}[\text{supp}(p_i, \mathbf{x}) + \text{supp}(p_{i+1}, \mathbf{x})] - \lambda \\
&= R^* + \frac{1}{2}[\text{supp}(p_i, \mathbf{x}) + \text{supp}(p_{i+1}, \mathbf{x})] - \lambda \quad (20)
\end{aligned}$$

where $\text{supp}(p_i, \mathbf{x})$, $\text{supp}(p_i, \mathbf{x})$ is the normalized support of p_i, p_{i+1} , defined in (3), and the regularization ‘bonus’ comes from the fact that d^* have one more leaf than d .

Because d^* is the optimal tree, we have $R^* \leq R(d, \mathbf{x}, \mathbf{y})$, which combined with (20) leads to (9). Therefore, for each child leaf pair p_k, p_{k+1} of a split, the sum of normalized supports of p_k, p_{k+1} should be no less than twice the regularization parameter, i.e., 2λ . \square

B.6 Proof of Theorem 3.7

PROOF. Let $d = (d'_{un}, \delta'_{un}, d'_{split}, \delta'_{split}, K', H')$ be the tree derived from d^* by deleting a pair of leaves $p_i \rightarrow \hat{y}_i^{(\text{leaf})}$, $p_{i+1} \rightarrow \hat{y}_{i+1}^{(\text{leaf})}$ and adding the their parent leaf $p_j \rightarrow \hat{y}_j^{(\text{leaf})}$. The discrepancy between d^* and d is the discrepancy between p_i, p_{i+1} and p_j : $\ell(d, \mathbf{x}, \mathbf{y}) - \ell(d^*, \mathbf{x}, \mathbf{y}) = a_i$, where a_i is defined in (10). Therefore,

$$\begin{aligned}
R(d, \mathbf{x}, \mathbf{y}) &= \ell(d, \mathbf{x}, \mathbf{y}) + \lambda(K - 1) = \ell(d^*, \mathbf{x}, \mathbf{y}) + a_i + \lambda(K - 1) \\
&= R(d^*, \mathbf{x}, \mathbf{y}) + a_i - \lambda = R^* + a_i - \lambda.
\end{aligned}$$

This combined with $R^* \leq R(d, \mathbf{x}, \mathbf{y})$ leads to $\lambda \leq a_i$. \square

B.7 Proof of Theorem 3.8

PROOF. Let $d = (d'_{un}, \delta'_{un}, d'_{split}, \delta'_{split}, K', H')$ be the tree derived from d^* by deleting a pair of leaves $p_i \rightarrow \hat{y}_i^{(\text{leaf})}$, $p_{i+1} \rightarrow \hat{y}_{i+1}^{(\text{leaf})}$ and adding the their parent leaf $p_j \rightarrow \hat{y}_j^{(\text{leaf})}$. The discrepancy between d^* and d is the discrepancy between p_i, p_{i+1} and p_j : $\ell(d, \mathbf{x}, \mathbf{y}) - \ell(d^*, \mathbf{x}, \mathbf{y}) = a_i$, where we defined a_i in (10). According to Theorem 3.7, $\lambda \leq a_i$ and

$$\begin{aligned}
\lambda &\leq \frac{1}{N} \sum_{n=1}^N \{ \text{cap}(x_n, p_i) \wedge \mathbb{1}[\hat{y}_i^{(\text{leaf})} = y_n] \\
&\quad + \text{cap}(x_n, p_{i+1}) \wedge \mathbb{1}[\hat{y}_{i+1}^{(\text{leaf})} = y_n] \\
&\quad - \text{cap}(x_n, p_j) \wedge \mathbb{1}[\hat{y}_j^{(\text{leaf})} = y_n] \}. \quad (19)
\end{aligned}$$

For any leaf j and its two child leaves $i, i + 1$, we always have

$$\begin{aligned}
\sum_{n=1}^N \text{cap}(x_n, p_i) \wedge \mathbb{1}[\hat{y}_i^{(\text{leaf})} = y_n] &\leq \sum_{n=1}^N \text{cap}(x_n, p_j) \wedge \mathbb{1}[\hat{y}_j^{(\text{leaf})} = y_n], \\
\sum_{n=1}^N \text{cap}(x_n, p_{i+1}) \wedge \mathbb{1}[\hat{y}_{i+1}^{(\text{leaf})} = y_n] &\leq \sum_{n=1}^N \text{cap}(x_n, p_j) \wedge \mathbb{1}[\hat{y}_j^{(\text{leaf})} = y_n]
\end{aligned}$$

which indicates that $a_i \leq \frac{1}{N} \sum_{n=1}^N \text{cap}(x_n, p_i) \wedge \mathbb{1}[\hat{y}_i^{(\text{leaf})} = y_n]$ and $a_i \leq \frac{1}{N} \sum_{n=1}^N \text{cap}(x_n, p_{i+1}) \wedge \mathbb{1}[\hat{y}_{i+1}^{(\text{leaf})} = y_n]$. Therefore,

$$\begin{aligned}
\lambda &\leq \frac{1}{N} \sum_{n=1}^N \text{cap}(x_n, p_i) \wedge \mathbb{1}[\hat{y}_i^{(\text{leaf})} = y_n], \\
\lambda &\leq \frac{1}{N} \sum_{n=1}^N \text{cap}(x_n, p_{i+1}) \wedge \mathbb{1}[\hat{y}_{i+1}^{(\text{leaf})} = y_n].
\end{aligned}$$

\square

B.8 Proof of Theorem 3.10

PROOF. By Corollary 3.4, $K \equiv \min(\lfloor 1/2\lambda \rfloor, 2^M)$ gives an upper bound on the number of leaves of any optimal tree. The algorithm begins by evaluating the empty tree, followed by M trees of depth $k = 1$, then $N_2 = \sum_{n_0=1}^1 \sum_{n_1=1}^{2^{n_0}} M \times \binom{2^{n_0}}{n_1} (M - 1)^{n_1}$ trees of depth $k = 2$. Before proceeding to length $k = 3$, we keep only $N_2 + C(M, 2) - P(M, 2)$ trees of depth $k = 2$, where N_k is defined in (1), $P(M, k)$ denotes the number of k -permutations of M and $C(M, k)$ denotes the number of k -combinations of M . Now, the number of length $k = 3$ prefixes we evaluate is $N_3 + C(M, 3) - P(M, 3)$. Propagating this forward gives (12) \square

B.9 Proof of Proposition 3.11

PROOF. Recall that the objective is $R(d, \mathbf{x}, \mathbf{y}) = \ell(d, \mathbf{x}, \mathbf{y}) + \lambda H$, where the misclassification error $\ell(d, \mathbf{x}, \mathbf{y})$ is given by

$$\ell(d, \mathbf{x}, \mathbf{y}) = \frac{1}{N} \sum_{n=1}^N \sum_{k=1}^K \text{cap}(x_n, p_k) \wedge \mathbb{1}[\hat{y}_k^{(\text{leaf})} \neq y_n].$$

Any particular tree uses a specific leaf, and therefore a single class label, to classify all points within a set of equivalent points. Thus, for a set of equivalent points u , the tree d correctly classifies either points that have the majority class label, or points that have the minority class label. It follows that d misclassifies a number of points in u at least as great as the number of points with the minority class label. To translate this into a lower bound on $\ell(d, \mathbf{x}, \mathbf{y})$, we first sum over all sets of equivalent points, and then for each such set, count differences between class labels and the minority class label of the set, instead of counting mistakes:

$$\begin{aligned}
\ell(d, \mathbf{x}, \mathbf{y}) &= \frac{1}{N} \sum_{u=1}^U \sum_{n=1}^N \sum_{k=1}^K \text{cap}(x_n, p_k) \wedge \mathbb{1}[\hat{y}_k^{(\text{leaf})} \neq y_n] \wedge \mathbb{1}[x_n \in e_u] \\
&\geq \frac{1}{N} \sum_{u=1}^U \sum_{n=1}^N \sum_{k=1}^K \text{cap}(x_n, p_k) \wedge \mathbb{1}[q_u = y_n] \wedge \mathbb{1}[x_n \in e_u].
\end{aligned}$$

Next, because every datum must be captured by a leaf in the tree d , $\sum_{k=1}^K \text{cap}(x_n, p_k) = 1$.

$$\ell(d, \mathbf{x}, \mathbf{y}) \geq \frac{1}{N} \sum_{u=1}^U \sum_{n=1}^N \mathbb{1}[x_n \in e_u] \mathbb{1}[y_n = q_u] = \sum_{u=1}^U \theta(e_u),$$

where the final equality applies the definition of $\theta(e_u)$ in (13). Therefore, $R(d, \mathbf{x}, \mathbf{y}) = \ell(d, \mathbf{x}, \mathbf{y}) + \lambda K \geq \sum_{u=1}^U \theta(e_u) + \lambda K$. \square