# PLANET: Massively Parallel Learning of Tree Ensembles with MapReduce

Biswanath Panda, Joshua S. Herbach, Sugato Basu, Roberto J. Bayardo

Google, Inc.

[bpanda, jsherbach, sugato]@google.com, bayardo@alum.mit.edu

## ABSTRACT

Classification and regression tree learning on massive datasets is a common data mining task at Google, yet many state of the art tree learning algorithms require training data to reside in memory on a single machine. While more scalable implementations of tree learning have been proposed, they typically require specialized parallel computing architectures. In contrast, the majority of Google's computing infrastructure is based on commodity hardware.

In this paper, we describe PLANET: a scalable distributed framework for learning tree models over large datasets. PLANET defines tree learning as a series of distributed computations, and implements each one using the *MapReduce* model of distributed computation. We show how this framework supports scalable construction of classification and regression trees, as well as ensembles of such models. We discuss the benefits and challenges of using a MapReduce compute cluster for tree learning, and demonstrate the scalability of this approach by applying it to a real world learning task from the domain of computational advertising.

## 1. INTRODUCTION

In this paper, we look at leveraging the MapReduce distributed computing framework for a complex data mining task of wide interest: learning ensembles of classification or regression trees. While there are other methods for parallel and distributed tree learning, building production-ready implementations remains complex and error-prone. With the wide and growing availability of MapReduce-capable compute infrastructures, it is natural to ask whether such infrastructures may be of use in parallelizing common data mining tasks such as tree learning. For many data mining operations, MapReduce may offer better scalability with vastly simplified deployment in a production setting.

MapReduce is a simple model for distributed computing that abstracts away many of the difficulties in parallelizing data management operations across a cluster of commodity machines. MapReduce reduces, if not eliminates, many com-

plexities such as data partitioning, scheduling tasks across many machines, handling machine failures, and performing inter-machine communication. These properties have motivated many technology companies to run MapReduce frameworks on their compute clusters for data analysis and other data management tasks. MapReduce has become in some sense an industry standard. For example, there are open source implementations such as Hadoop that can be run either in-house or on cloud computing services such as Amazon EC2.[1] Startups like Cloudera[2] offer software and services to simplify Hadoop deployment, and companies including Google, IBM and Yahoo! have granted several universities access to Hadoop clusters to further cluster computing research.[3]

Despite the growing popularity of MapReduce [12], its application to certain standard data mining and machine learning tasks remains poorly understood. In this paper we focus on one such task: tree learning. We believe that a tree learner capable of exploiting a MapReduce cluster can effectively address many scalability issues that arise in building tree models on massive datasets. Our choice of focusing on tree models is motivated primarily by their popularity. Tree models are used in many applications because they are interpretable, can model complex interactions, and can handle both ordered and unordered features. Recent studies have shown that tree models, when combined with ensemble techniques, provide excellent predictive performance across a wide variety of domains [8, 9].

This paper describes our experiences with developing and deploying a MapReduce based tree learner called PLANET, which stands for Parallel Learner for Assembling Numerous Ensemble Trees. The development of PLANET was motivated by a real application in sponsored search advertising in which massive clickstreams are processed to develop a predictor of user experience following the click of a sponsored search ad [30]. We show how PLANET can be scaled effectively to large datasets, describe experiments that highlight the performance characteristics of PLANET, and demonstrate the benefits of various optimizations that we implemented within the system. We show that while MapReduce is not a panacea, it still provides a powerful basis on which scalable tree learning can be implemented.

The rest of the paper is organized as follows. In Section 2 we describe the necessary background on which we build,

---

[1] http://aws.amazon.com/ec2/

[2] http://www.cloudera.com/

[3] For example, see http://www.youtube.com/watch?v=UBrDPRlplyo and http://www.nsf.gov/news/news_summ.jsp?cntn_id=111470

including the formal problem definitions of classification and regression. We also review the process of solving these problems through tree induction, and describe the MapReduce paradigm for distributed computation. As a prelude to a more detailed description of our approach, in Section 3 we provide an example of how tree induction proceeds in PLANET. This example describes the roles of each of the major components as well as their high level requirements. Section 4 provides a more formal algorithm description for the case of learning a single classification or regression tree, and Section 5 describes how PLANET can be generalized to produce ensembles of trees via boosting and/or bagging. In Section 6 we discuss several important details we had to address in our efforts to develop an efficient and production-ready deployment. We describe the performance of PLANET on our sponsored search derived clickstream dataset in Section 7. We review related work in Section 8 and conclude with a discussion of future work in Section 9.

## 2. PRELIMINARIES

Let $\mathcal{X} = \{X_1, X_2, \ldots X_N\}$ be a set of attributes with domains $\mathbb{D}_{X_1}, \mathbb{D}_{X_2}, \ldots \mathbb{D}_{X_N}$ respectively. Let $Y$ be an output with domain $\mathbb{D}_Y$. Consider a dataset $D^* = \{(\mathbf{x}_i, y_i) | \mathbf{x}_i \in \mathbb{D}_{X_1} \times \mathbb{D}_{X_2} \times \ldots \mathbb{D}_{X_N}, y_i \in \mathbb{D}_Y\}$ sampled from an unknown distribution, where the $i^{th}$ data vector $\mathbf{x}_i$ has an output $y_i$ associated with it. Given the dataset $D^*$, the goal in supervised learning is to learn a function (or *model*) $F : \mathbb{D}_{X_1} \times \mathbb{D}_{X_2} \times \ldots \mathbb{D}_{X_N} \to \mathbb{D}_Y$ that best approximates the true distribution of $D^*$. If $\mathbb{D}_Y$ is continuous, the learning problem is a regression problem; if $\mathbb{D}_Y$ is categorical, it is a classification problem.

Let $\mathcal{L}$ be a function that quantifies in some way the discrepancy between the function prediction $F(\mathbf{x}_i)$ on $\mathbf{x}_i$ and the actual output $y_i$. A model that minimizes the net loss $\sum_{(\mathbf{x}_i, y_i) \in D^*} \mathcal{L}(F(\mathbf{x}_i), y_i)$ on the *training set* $D^*$ may not generalize well (have low loss) when applied to future data [32]. Generalization is attained through controlling model complexity by various methods, e.g., pruning and ensemble learning for tree models [5]. The learned model is evaluated by measuring its net loss when applied to a holdout data set.

### 2.1 Tree Models

Classification and regression trees are one of the oldest and most popular data mining models [13]. Tree models represent $F$ by recursively partitioning the data space $\mathbb{D}_{X_1} \times \mathbb{D}_{X_2} \times \ldots \mathbb{D}_{X_N}$ into non-overlapping regions, with a simple model in each region.

Figure 1 shows an example tree model. Non-leaf nodes in the tree define region boundaries in the data space. Each region boundary is represented as a predicate on an attribute in $\mathcal{X}$. If the attribute is ordered, the predicate is of the form $X < v$, $v \in \mathbb{D}_X$ (e.g., Node A in Figure 1). Unordered attributes have predicates of the form $X \in \{v_1, v_2, \ldots v_k\}$, $v_1 \in \mathbb{D}_X, v_2 \in \mathbb{D}_X, \ldots v_k \in \mathbb{D}_X$, (e.g., Node B in Figure 1). The path from the root to a leaf node in the tree defines a region. Leaf nodes (e.g., the left child of A in Figure 1), contain a region prediction which in most cases is a constant value or some simple function. To make predictions on an unknown $\mathbf{x}$, the tree is traversed to find the region containing $\mathbf{x}$. The region containing $\mathbf{x}$ is the path from the root to a leaf in the tree along which all non-leaf predicates are true when evaluated on $\mathbf{x}$. The prediction given by this leaf is used as the value for $F(\mathbf{x})$.

---

**Algorithm 1** InMemoryBuildNode

**Require:** Node $n$, Data $D \subseteq D^*$
1: $(n \to \text{split}, D_L, D_R) = \text{FindBestSplit}(D)$
2: **if** StoppingCriteria$(D_L)$ **then**
3:    $n \to \text{left\_prediction} = \text{FindPrediction}(D_L)$
4: **else**
5:    InMemoryBuildNode$(n \to \text{left}, D_L)$
6: **if** StoppingCriteria$(D_R)$ **then**
7:    $n \to \text{right\_prediction} = \text{FindPrediction}(D_R)$
8: **else**
9:    InMemoryBuildNode$(n \to \text{right}, D_R)$

---

In our example tree model, predicate evaluations at non-leaf nodes have only two outcomes, leading to binary splits. While tree models can have non-binary splits, for the sake of simplicity we will focus on binary splits only for the remainder of this paper. All our techniques also apply to tree algorithms with non-binary splits with straightforward modifications.

Tree models are popular because they are interpretable, capable of modeling complex classification and regression tasks, and handle both ordered and categorical domains. Recent work by Caruana et al. [9] has also shown that tree models, when combined with ensemble learning methods like bagging [4], boosting [14], and forests [5], outperform many other popular learning methods in terms of prediction accuracy. A thorough discussion of tree models and different ensemble methods is beyond the scope of this paper — see [29] for a good review.

### 2.2 Learning Tree Models

Previous work on learning tree models is extensive. For a given training dataset $D^*$, finding the optimal tree is known to be NP-Hard; thus most algorithms use a greedy top-down approach to construct the tree (Algorithm 1) [13]. At the root of the tree, the entire training dataset $D^*$ is examined to find the *best* split predicate for the root. The dataset is then partitioned along the split predicate and the process is repeated recursively on the partitions to build the child nodes.

Finding the best split predicate for a node (Line 1) is the most important step in the greedy learning algorithm, and has been the subject of much of the research in tree learning. Numerous techniques have been proposed for finding the right split at a node, depending on the particular learning problem. The main idea is to reduce the *impurity* ($I$) in a node. Loosely defined, the impurity at a node is a measure of the dissimilarity in the $Y$ values of the training records $D$ that are input to the node. The general strategy is to pick a predicate that maximizes $I(D) - (I(D_L) + I(D_R))$, where $D_L$ and $D_R$ are the datasets obtained after partitioning $D$ on the chosen predicate. At each step the algorithm greedily partitions the data space to progressively reduce region impurity. The process continues until all $Y$ values in the input dataset $D$ to a node are the same, at which point the algorithm has isolated a pure region (Lines 2-3 and 6-7). Some algorithms do not continue splitting until regions are completely pure, and instead stop once the number of records in $D$ falls below a predefined threshold.

Popular impurity measures that have been proposed are derived from measures such as entropy, Gini index, and variance [29], to name only a few. PLANET uses an impurity

measure based on variance ($Var$) to evaluate the quality of a split. The higher the variance in the $Y$ values of a node, the greater the node's impurity. Further details on the split criteria are discussed in Section 2.3. While we focus concretely on variance as our split criteria for the remainder of this presentation, as long as a split metric can be computed on subsets of the training data and later aggregated, PLANET can be easily extended to support it.

### 2.2.1 Scalability Challenge

The greedy tree induction algorithm we have described is simple and works well in practice. However, it does not scale well to large training datasets. FindBestSplit requires a full scan of the node's input data, which can be large at higher levels of the tree. Large inputs that do not fit in main memory become a bottleneck because of the cost of scanning data from secondary storage. Even at lower levels of the tree where a node's input dataset $D$ is typically much smaller than $D^*$, loading $D$ into memory still requires reading and writing partitions of $D^*$ to secondary storage multiple times.

Previous work has looked at problem of building tree models from datasets which are too large to fit completely in main memory. Some of the known algorithms are disk-based approaches that use clever techniques to optimize the number of reads and writes to secondary storage during tree construction (e.g., [26]). Other algorithms scan the training data in parallel using specialized parallel architectures (e.g., [3]). We defer a detailed discussion of these approaches and how they compare to PLANET to Section 8. As we will show in Section 8, some of the ideas used in PLANET have been proposed in the past; however, we are not aware of any efforts to build massively parallel tree models on commodity hardware using the MapReduce framework.

Post-pruning learned trees to prevent overfitting is also a well studied problem. However, with ensemble models (Section 5), post pruning is not always needed. Since PLANET is primarily used for building ensemble models, we do not discuss post pruning in this paper.

## 2.3 Regression Trees

Regression trees are a special case of tree models where the output attribute $Y$ is continuous [5]. We focus primarily on regression trees within this presentation because most of our use cases require predictions on continuous outputs. Note that any regression tree learner also supports binary (0-1) classification tasks by modeling them as instances of logistic regression. The core operations of regression tree learning in Algorithm 1 are implemented as follows:

**FindBestSplit**($D$): In a regression tree, $D$ is split using the predicate that results in the largest reduction in variance. Let $Var(D)$ be the variance of the output attribute $Y$ measured over all records in $D$. At each step the tree learning algorithm picks a split which maximizes

$$|D| \times Var(D) - (|D_L| \times Var(D_L) + |D_R| \times Var(D_R)), \quad (1)$$

where $D_L \subset D$ and $D_R \subset D$ are the training records in the left and right subtree after splitting $D$ by a predicate.

Regression trees use the following policy to determine the set of predicates whose split quality will be evaluated:

- For ordered domains, split predicates are of the form $X_i < v$, for some $v \in \mathbb{D}_{X_i}$. To find the best split, $D$ is sorted along $X_i$, and a split point is considered between each adjacent pair of values for $X_i$ in the sorted list.

- For unordered domains, split predicates are of the form $X_i \in \{v_1, v_2, \ldots v_k\}$, where $\{v_1, v_2, \ldots v_k\} \in \mathcal{P}(\mathbb{D}_{X_i})$, the power set of $\mathbb{D}_{X_i}$. Breiman [6] presents an algorithm for finding the best split predicate for a categorical attribute without evaluating all possible subsets of $\mathbb{D}_{X_i}$. The algorithm is based on the observation that the optimal split predicate is a subsequence in the list of values for $X_i$ sorted by the average $Y$ value.

**StoppingCriteria**($\mathbf{D}$): A node in the tree is not expanded if the number of records in $D$ falls below a threshold. Alternatively, the user can also specify the maximum depth to which a tree should be built.

**FindPrediction**($\mathbf{D}$): The prediction at a leaf is simply the average of the all the $Y$ values in $D$.

## 2.4 MapReduce

PLANET uses MapReduce [12] to distribute and scale tree induction to very large datasets. MapReduce provides a framework for performing a two-phase distributed computation on large datasets, which in our case is the training dataset $D^*$. In the *Map* phase, the system partitions $D^*$ into a set of disjoint units which are assigned to workers, known as mappers. In parallel, each mapper scans through its assigned data and applies a user-specified map function to each record. The output of the user's map function is a set of $\langle key, value \rangle$ pairs which are collected for MapReduce's *Reduce* phase. In the reduce phase, the key-value pairs are grouped by key and are distributed to a series of workers, called reducers. Each reducer then applies a user-specified reduce function to all the values for a key and outputs a final value for the key. The collection of final values from all of the reducers is the final output of MapReduce.

## 3. EXAMPLE

The PLANET framework breaks up the process of constructing a tree model into a set of MapReduce tasks. Dependencies exist between the different tasks, and PLANET uses clever scheduling methods to efficiently execute and manage them. Before delving into the technical details of the framework, we begin with a detailed example of how tree induction proceeds in PLANET.

The example introduces the different components in PLANET, describes their roles, and provides a high level overview of the entire system. To keep the example simple we only discuss the construction of a single tree. The method extends naturally to ensembles of trees, as we discuss in Section 5.

**Example setup:** Let us assume that we have a training dataset $D^*$ with 100 records. Further assume that tree induction stops once the number of training records at a node falls below 10. Let the tree in Figure 1 be the model that will be learned if we ran Algorithm 1 on a machine with sufficient memory. Our goal in this example is to demonstrate how PLANET constructs the tree in Figure 1 when there is a memory constraint limiting Algorithm 1 to operating on inputs of size 25 records or less.
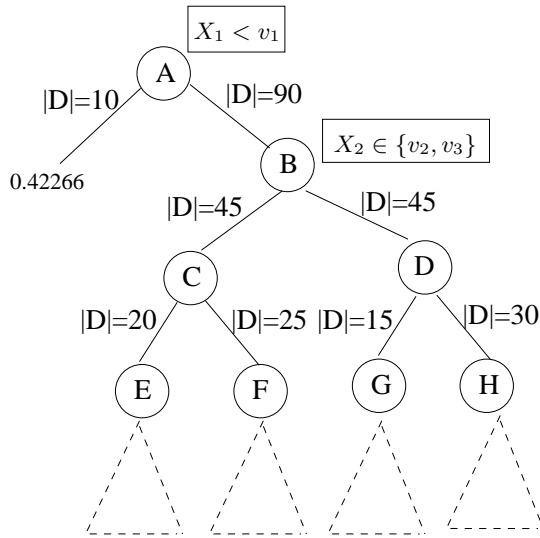
**Figure 1: Example Tree. Note that the labels on the nodes (in boxes) are the split predicates, while the labels on the edges are the sizes of the dataset in each branch (|D| denotes the dataset size in that branch in this figure).**

## 3.1 Components

At the heart of PLANET is the *Controller*, a single machine that initiates, schedules and controls the entire tree induction process. The Controller has access to a compute cluster on which it schedules MapReduce jobs. In order to control and coordinate tree construction, the Controller maintains the following:

- *ModelFile* (M): The Controller constructs a tree using a set of MapReduce jobs, each of which builds different parts of the tree. At any point, the model file contains the entire tree constructed so far.

Given the ModelFile (M), the Controller determines the nodes at which split predicates can be computed. In the example of Figure 1, if M has nodes A and B, then the Controller can compute splits for C and D. This information is stored in two queues.

- *MapReduceQueue* (MRQ): This queue contains nodes for which $D$ is too large to fit in memory (i.e. $> 25$ in our example).

- *InMemoryQueue* (InMemQ): This queue contains nodes for which $D$ fits in memory (i.e $\leq 25$ in our example).

As tree induction proceeds, the Controller dequeues nodes off MRQ and InMemQ and schedules MapReduce jobs to find split predicates at the nodes. Once a MapReduce job completes, the Controller updates M with the nodes and their split predicates, and then updates MRQ and InMemQ with new nodes at which split predicates can be computed. Each MapReduce job takes as input a set of nodes ($N$), the training data set ($D^*$), and the current state of the model (M). The Controller schedules two types of MapReduce jobs.

- Nodes in MRQ are processed using *MR_ExpandNodes*, which for a given set of nodes $N$ computes a candidate set of good split predicates for each node in $N$.

- Nodes in InMemQ are processed using *MR_InMemory*. Recall that nodes in InMemQ have input data sets $D$ that are small enough to fit in memory. Therefore, given a set of nodes $N$, MR_InMemory completes tree induction at nodes in $N$ using Algorithm 1.

We defer details of the MapReduce jobs to the next section. In the remainder of this section, we will tie the above components together and walk through the example.

## 3.2 Walkthrough

When tree induction begins, M, MRQ, and InMemQ are all empty. The only node the Controller can expand is the root (A). Finding the split for A requires a scan of the entire training dataset of 100 ($\geq 25$) records. Since this set is too large to fit in memory, A is pushed onto MRQ and InMemQ stays empty.

After initialization the Controller dequeues A from MRQ and schedules a job MR_ExpandNodes({A}, M, $D^*$). This job computes a set of good splits for node A along with some additional information about each split. Specifically, for each split we compute (1) the quality of the split (i.e., the reduction in impurity), (2) the predictions in the left and right branches, and (3) the number of training records in the left and right branches.

The split information computed by MR_ExpandNodes gets sent back to the Controller, which selects the best split for node A. In this example, the best split has 10 records in the left branch and 90 records in the right. The selected split information for node A is then added into the ModelFile. The Controller next updates the queues with new nodes at which split predicates can be computed. The left branch of A has 10 records. This matches the stopping criteria and hence no new nodes are added for this branch. For the right branch with 90 records ($\geq 25$), node B can be expanded and is pushed onto MRQ.

Tree induction continues by dequeuing node B, and scheduling MR_ExpandNodes({B}, M, $D^*$). Note that for expanding node B we only need the records that went down the right subtree of A, but to minimize book keeping, PLANET passes the entire training dataset to the MapReduce. As we describe in 4.3, MR_ExpandNodes uses the current state of the ModelFile to determine the subset of $D^*$ that will be input to B.

Once the Controller has received the results for the MapReduce on node B and updated M with the split for B, it can now expand both C and D. Both of these nodes get 45 records as input and are therefore pushed on to MRQ. The Controller can now schedule a single MR_ExpandNodes({C, D}, M, $D^*$) job to find the best splits for both nodes C and D. Note that by expanding C and D in a single step, PLANET expands trees breadth first as opposed to the depth first process used by the in-memory Algorithm 1.

Once the Controller has the obtained the splits for C and D, it can schedule jobs to expand nodes E, F, G, and H. Of these, H uses 30 records, which still cannot fit in memory, and hence gets added to MRQ. The input sets to E, F, G are small enough to fit into memory and hence tree induction at these nodes can be completed in-memory. The Controller pushes these nodes into the InMemQueue.

The Controller next schedules two MapReduce jobs simultaneously. MR_InMemory({E,F,G}, M, $D^*$) completes tree induction at nodes E, F, and G since the input datasets to these nodes are small. MR_ExpandNodes({H}, M, $D^*$)

---

**Algorithm 2** MR_ExpandNodes::Map

---

**Require:** NodeSet $N$, ModelFile M, Training record $(\mathbf{x}, y) \in D^*$
1: $n = \text{TraverseTree}(M, \mathbf{x})$
2: **if** $n \in N$ **then**
3:     agg_tup$_n \leftarrow y$
4:     **for all** $X \in \mathcal{X}$ **do**
5:       $v = \text{Value on } X \text{ in } \mathbf{x}$
6:       **if** $X$ is ordered **then**
7:         **for all** Split point $s$ of $X$ s.t. $s \leqslant v$ **do**
8:           $T_{n,X}[s] \leftarrow y$
9:       **else**
10:        $T_{n,X}[v] \leftarrow y$

---

---

**Algorithm 3** MR_ExpandNodes::Map_Finalize

---

**Require:** NodeSet $N$
1: **for all** $n \in N$ **do**
2:     Output to all reducers(agg_tup$_n$)
3:     **for all** $X \in \mathcal{X}$ **do**
4:       **if** $X$ is ordered **then**
5:         **for all** Split point $s$ of $X$ **do**
6:           Output$((n, X, s), T_{n,X}[s])$
7:       **else**
8:         **for all** $v \in T_{n,X}$ **do**
9:           Output$((n, X), (v, T_{n,X}[v]))$

---

computes good splits for H. Once the InMemory job returns, tree induction for the subtrees rooted at E, F, and G is complete. The Controller updates MRQ and InMemQ with the children of node H and continues tree induction. PLANET aggressively tries to maximize the number of nodes at which split predicates can be computed in parallel, and schedules multiple MapReduce jobs simultaneously.

## 4. TECHNICAL DETAILS

In this section, we discuss the technical details of PLA-NET's major components — the two critical MapReduces that handle splitting nodes and growing subtrees, and the Controller that manages the entire tree induction process.

### 4.1 MR_ExpandNodes: Expanding a Single Node

MR_ExpandNodes is the component that allows PLANET to train on datasets too large to fit in memory. Given a set of nodes ($N$), the training dataset ($D^*$), and the current model ($M$), this MapReduce job computes a set of good splits for each node in $N$.

**Map Phase:** The training dataset $D^*$ is partitioned across a set of mappers. Each mapper loads into memory the current model (M) and the input nodes $N$. Note that the union of the input datasets to all nodes in $N$ need not be equal to $D^*$. However, every MapReduce job scans the entire training data set applying a Map function to every training record. We will discuss this design decision in Section 4.3. Pseudocode describing the algorithms that are executed by each mapper appear in Algorithms 2 and 3. Given a training record $(\mathbf{x}, y)$, a mapper first determines if the record is part of the input dataset for any node in $N$ by traversing the current model M with $(\mathbf{x}, y)$ (Line 1, Alg. 2). Once the input set to a node is determined, the next step is to evaluate

possible splits for the node, and select the best one.

Recall from Section 2.3 the method for finding the best split for a node $n$. For an ordered attribute $X$, Equation 1 is computed between every adjacent pair of values for the attribute that appear in the node's input dataset $D$. Performing this operation in a distributed setting would require us to sort $D^*$ along each ordered attribute and write out the results to secondary storage. These sorted records would then have to be partitioned carefully across mappers, keeping track of the range of values on each mapper. Distributed algorithms implementing such approaches are complex and end up using additional storage or network resources. PLA-NET makes a tradeoff between finding the perfect split for an ordered attribute and simple data partitioning. Splits are not evaluated between every pair of values of an attribute. Rather, prior to tree induction we run a MapReduce on $D^*$ and compute approximate equidepth histograms for every ordered attribute [25]. When computing splits on an ordered attribute, a single split point is considered from every histogram bucket of the attribute.

On startup, each mapper loads the set of split points to be considered for each ordered attribute. For each node $n \in N$ and attribute $X$, the mapper maintains a table $T_{n,X}$ of key-value pairs. Keys for the table are the split points to be considered for $X$ and the values are tuples ($agg\_tup$) of the form $\{\sum y, \sum y^2, \sum 1\}$. For a particular split point $s \in \mathbb{D}_X$ being considered for node $n$, the tuple $T_{n,X}[s]$ contains: (1) the sum of $Y$ values for training records $(\mathbf{x}, y)$ that are input to $n$ and have values for $X$ that are less than $s$, (2) the sum of squares of these values, and (3) the number of training records that are input to $n$ and have values of $X$ less than $s$. Mappers scan subsets of $D^*$ and compute agg_tups for all split points being considered for each node in $N$ (Lines 7, 8 in Alg. 2). After processing all its input records, each mapper outputs keys of the form $n, X, s$ and the corresponding $T_{n,X}[s]$ as values (Line 6, Alg. 3). As we show later, a reduce function will aggregate the agg_tups with the same key to compute the quality of the split $X < s$ for node $n$.

For computing splits on an unordered attribute $X$, Section 2.3 proposed computing Equation 1 for every subsequence of unique values of $X$ sorted by the average $Y$. Each mapper performs this computation by maintaining a table $T_{n,X}$ of key, agg_tup pairs as described before. However, in this case keys correspond to unique values of $X$ seen in the input records to node $n$. $T_{n,X}[v]$ maintains the same aggregate statistics as described earlier for all training records that are input to $n$ and have an $X$ value of $v$ (Line 10, Alg. 2). After processing all input data, the mappers output keys of the form $n, X$ and value $\langle v, T_{n,X}[v] \rangle$ (Line 9, Alg. 3). Note the difference in key-value pairs output for ordered and unordered attributes. Quality of a split on an ordered attribute can be computed independently of other splits on that attribute, hence the split point $s$ is part of the key. To run Breiman's algorithm, all values of an unordered attribute need to be sorted by average $Y$ value. Hence, the value $v$ of an attribute is not part of the key. A single reducer processes and sorts all the values of the attribute to compute the best split on the attribute.

In addition to the above outputs, each mapper also maintains agg_tup$_n$ for each node $n \in N$ (Line 3, Alg. 2) and outputs them to all reducers (Line 2, Alg. 3). These tuples are computed over all input records to their respective nodes, and help reducers in computing split qualities.

**Algorithm 4** MR_ExpandNodes::Reduce

---
**Require:** Key $k$,Value Set $V$
 1: **if** $k == n$ **then**
 2:  {Aggregate agg_tup$_n$'s from mappers}
 3:  agg_tup$_n$ = Aggregate($V$)
 4: **else if** $k == n, X, s$ **then**
 5:  {Split on ordered attribute}
 6:  agg_tup$_{left}$ = Aggregate($V$)
 7:  agg_tup$_{right}$ = agg_tup$_n$ - agg_tup$_{left}$
 8:  UpdateBestSplit($S[n]$,$X$,$s$,agg_tup$_{left}$, agg_tup$_{right}$)
 9: **else if** $k == n, X$ **then**
10:  {Split on unordered attribute}
11:  **for all** $v$,agg_tup $\in$ V **do**
12:   $T[v] \leftarrow$ agg_tup
13:  UpdateBestSplit($S[n]$,BreimanSplit($X$,$T$,agg_tup$_n$))

---

**Reduce Phase:** The reduce phase, which works on the outputs from the mappers, performs aggregations and computes the quality of each split being considered for nodes in $N$. Each reducer maintains a table $S$ indexed by nodes. $S[n]$ contains the best split seen by the reducer for node $n$. The pseudocode executed on each reducer is outlined in Algorithm 4. A reducer processes three types of keys. The first is of the form $n$ with a value list $V$ of the all agg_tup$_n$ tuples output by the mappers. These agg_tups are aggregated to get a single agg_tup$_n$ with the $\{\sum y, \sum y^2, \sum 1\}$ values for all input records to node $n$ (Line 3, Alg. 4). Reducers process keys in sorted order so that they process all keys of type $n$ first. The other types of keys that a reducer processes belong to ordered and unordered attributes. The keys corresponding to unordered attributes are of the form $n, X$. Here the set $V$ associated with each key is a set of pairs consisting of an unordered attribute value $v$ and an agg_tup. For each $v$ the agg_tups are aggregated to get $\{\sum y, \sum y^2, \sum 1\}$ over all input records to $n$ where the value of $X$ is $v$. Once aggregated, Breiman's algorithm is used to find the optimal split for $X$, and $S[n]$ is updated if the resulting split is better than any previous split for $n$ (Lines 11-13, Alg 4). For ordered attributes, keys are of the form $n, X, s$ and $V$ is again a list of agg_tups. Aggregating these into agg_tup$_{left}$ gives the $\{\sum y, \sum y^2, \sum 1\}$ values for all records input to $n$ that fall in the left branch of $X < s$ (Line 6, Alg. 4). Using agg_tup$_n$ and agg_tup$_{left}$ it is straightforward to compute the $Var$ based quality of the split $X < s$. If this split $X < s$ is better than the best split seen by the reducer for $n$ so far, then $S[n]$ is updated to the current split (Lines 7-8, Alg. 4).

Finally, each reducer outputs the best split $S[n]$ that it has seen for each node. In addition to the split quality and predicate, it also outputs the average $Y$ value, and number of the training records in the left and right branches of the split. The Controller takes the splits produced by all the reducers and finds the best split for each node in $N$, then updates the ModelFile M with this information. The Controller updates the queues with the child nodes that should be expanded using information about the number of training records in each branch.

## 4.2 MR_InMemory: In Memory Tree Induction

As tree induction progresses, the size of the input dataset for many nodes becomes small enough to fit in memory.

**Algorithm 5** UpdateQueues

---
**Require:** DataSetSize $|D|$, Node $n$
 1: **if** not StoppingCriteria($|D|$) **then**
 2:  **if** $|D| <$ in_memory_threshold **then**
 3:   InMemQ.append($n$)
 4:  **else**
 5:   MRQ.append($n$)

---

**Algorithm 6** Schedule_MR_ExpandNode

---
**Require:** NodeSet $N$,Current Model M
 1: CandidateGoodSplits = MR_ExpandNodes($N$,M,$D^*$)
 2: **for all** $n \in N$ **do**
 3:  $n \to$split,$n \to$l_pred, $n \to$r_pred,$|D_L|$,$|D_R| =$
      FindBestSplit($n$, CandidateGoodSplits)
 4:  UpdateQueues($|D_L|$,$n \to$left)
 5:  UpdateQueues($|D_R|$,$n \to$right)
 6: jobs_running - -

---

At any such point, rather than continuing tree induction using MR_ExpandNodes, the Controller completes tree induction in-memory using a different MapReduce job called MR_InMemory. Like MR_ExpandNodes, MR_InMemory partitions $D^*$ across a set of mappers. The map function processes a training record $(\mathbf{x}, y)$ and traverses the tree in M, to see if the $(\mathbf{x}, y)$ is input to some node $n \in N$. If such a node is found then the map function outputs the node $n$ as the key and $(\mathbf{x}, y)$ as the value. The reduce function receives as input a node $n$ (as key) and the set of training records that are input to the node (as values). The reducer loads the training records for $n$ into memory and completes subtree construction at $n$ using Algorithm 1.

## 4.3 Controller Design

The example in Section 3 provides the intuition behind functionality of the Controller. Here we provide a more detailed look at its roles and implementation.

The main Controller thread (Algorithm 8) schedules jobs off of its queues until the queues are empty and none of the jobs it schedules remain running. Scheduled MapReduce jobs are launched in separate threads so that the Controller can send out multiple jobs in parallel. When a MR_ExpandNodes job returns, the queues are updated with the new nodes that can now be expanded (Algorithm 6). Note that when MR_InMemory finishes running on a set of nodes $N$ (Algorithm 7), no updates are made to the queues because tree induction at nodes in $N$ is complete.

While the overall architecture of the Controller is fairly straightforward, we would like to highlight a few important design decisions. First, in our example in Section 3, recall that the Controller always removed all existing nodes from MRQ and InMemQ and scheduled MapReduce jobs. Therefore, it may seem that the Controller need not maintain queues and can schedule subsequent MapReduce jobs directly after processing the output of a MapReduce job.

**Algorithm 7** Schedule_MR_InMemory

---
**Require:** NodeSet $N$,Current Model M
 1: MR_InMemory($N$,M,$D$)
 2: jobs_running - -

---

---

**Algorithm 8** MainControllerThread

---
**Require:** Model M = ∅, MRQ=∅, InMemQ=∅
1: MRQ.append(root)
2: **while** true **do**
3:   **while** MRQ not empty **do**
4:     **if** TryReserveClusterResources **then**
5:       jobs_running ++
6:       NewThread(ScheduleMR_ExpandNode(⊆MRQ,M))
7:   **while** InMemQ not empty **do**
8:     **if** TryReserveClusterResources **then**
9:       jobs_running ++
10:       NewThread(ScheduleMR_InMemory(⊆InMemQ,M))
11:   **if** jobs_running==0 && MRQ empty && InMemQ empty **then**
12:     Exit

---

However, in practice this is not always possible. The memory limitations on a machine and the number of available machines on the cluster often prevent the Controller from scheduling MapReduce jobs for all nodes on a queue at once.

Second, when scheduling a set of nodes, recall that the Controller does not determine the set of input records required by the nodes. Instead, it simply sends the entire training dataset $D^*$ to every job. If the input to the set of nodes being expanded by a node is much smaller than $D^*$, then this implementation results in the Controller sending much unnecessary input for processing. On the other hand, this design keeps the overall system simple. In order to avoiding sending unnecessary input, the Controller would need to write out the input training records for each node to storage. This in turn would require additional bookkeeping for the Controller when operating normally, and would further complicate important systems like the checkpointing mechanism (Section 6.3) and ensemble creation (Section 5). The amount of unnecessary information sent by our implementation is also mitigated by breadth-first tree construction. If we can expand all nodes at level $i + 1$ in one MapReduce job, then every training record is part of the input to some node that is being expanded. Finally, MapReduce frameworks are already optimized for scanning data efficiently in a distributed fashion – the additional cost of reading in a larger dataset can be mitigated by adding more mappers, if necessary.

## 5. LEARNING ENSEMBLES

Until now we have described how the PLANET framework builds a single tree. Ensemble-based tree models have better predictive power when compared to single tree models [8, 9]. Bagging [4] and boosting [15] are the two most popular tree ensemble learning methods. In this section we show how PLANET supports the construction of tree ensembles through these two techniques.

Boosting is an ensemble learning technique that uses a weighted combination of weak learners to form a highly accurate predictive model [14]. Our current boosting implementation uses the GEM algorithm proposed by Friedman [15]. In the GEM algorithm, every weak learner is a shallow tree (depth ≈ 2 or 3). Model construction proceeds as follows: assume $k - 1$ weak learners (shallow trees) have been added to the model. Let $F_{k-1}$ be the boosted model of those trees. Tree $k$ is trained on a sample of $D^*$ and residual

predictions ($z$). For a given training record $(\mathbf{x}, y)$, the residual prediction for tree $k$ is $z = y - F_{k-1}(\mathbf{x})$ for a regression problem, and $z = y - \frac{1}{1+exp(-F_{k-1}(\mathbf{x}))}$ for a classification problem. The boosting process is initialized by setting $F_0$ as some aggregate defined over the $Y$ values in the training dataset. Abstracting out the details, we need three main features in our framework to build boosted models.

- **Building multiple trees:** Extending the Controller to build multiple trees is straightforward. Since the Controller manages tree induction by reducing the process to repeated node expansion, the only change necessary for constructing a boosted model is to push the root node for tree $k$ onto the MR after tree $k - 1$ is completed.

- **Residual computation: Training trees on residuals is simple since the current model is sent to every Map Reduce job in full.** If the mapper decides to use a training record as input to a node, it can compute the current model's prediction, and hence the residual.

- **Sampling**: Each tree is built on a sample of $D^*$. Dappers compute a hash of a training record's id and the tree id. Records hashing into a particular range are used for constructing the tree. This hash-based sampling guarantees that the same sample will be used for all nodes in a tree, but different samples of $D^*$ will be used for different trees.

Building an ensemble model using bagging involves learning multiple trees over independent samples of the training data. Predictions from each tree in the model are computed and averaged to compute the final model prediction. PLANET supports bagging as follows: when tree induction begins at the root, nodes of all trees in the bagged model are pushed onto the MRQ. The Controller then continues tree induction over dataset samples as already described. In this scenario, at any point in time the queues will contain nodes belonging to many different trees instead of a single tree, thereby allowing the Controller to exploit greater parallelism.

The bagging algorithm proposed by Breiman expects each tree to be built on a sample of $D^*$ generated with replacement; however, our framework only supports sampling without replacement at this time. We are still exploring efficient techniques to do sampling with replacement in a distributed setting, and comparing how bagged models generated by sampling with and without replacement differ in practice.

## 6. ENGINEERING ISSUES

In developing a production-capable deployment of PLANET, we encountered several unanticipated challenges. First, because MapReduce was not intended to be used for highly iterative procedures like tree learning, we found that MapReduce start up and tear down costs were primary performance bottlenecks. Second, the cost of traversing models in order to determine split points in parallel turned out to be higher than we expected. Finally, even though MapReduce offers graceful handling of failures within a specific MapReduce computation, since our computation spans multiple MapReduce phases, dealing with shared and unreliable commodity resources remained an issue which we had to address. We discuss our solutions to each of these issues within this section.
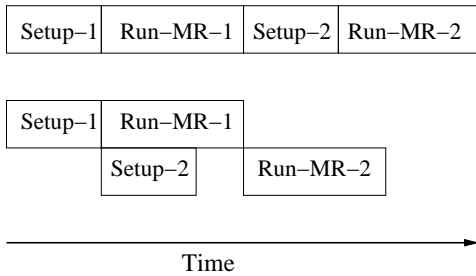
| Setup–1 | Run–MR–1 | Setup–2 | Run–MR–2 |
|---|---|---|---|

| Setup–1 | Run–MR–1 | |
|---|---|---|
| | Setup–2 | Run–MR–2 |

Time

**Figure 2: Forward Scheduling**

## 6.1 Forward Scheduling

Immediately after our initial attempt at deploying PLA-NET on a live map reduce cluster, we noticed that an inordinate amount of time was spent in setting up and tearing down MapReduce jobs. Fixing latency due to tear down time was a simple change to the logic in Algorithms 6 and 8. Instead of waiting for a MapReduce job to finish running on the cluster, the Controller ran a thread which would periodically check for the MapReduce's output files. Once the output files were available, the thread would load them and run the FindBestSplit and UpdateQueues logic described in Algorithm 6.

Addressing the latency caused by job set up was a more interesting challenge. Set up costs include time spent allocating machines for the job, launching a master to monitor the MapReduce job, and preparing and partitioning the input data for the MapReduce. To get around this problem we implemented a simple trick of forward scheduling MapReduce jobs. Figure 2 illustrates the basic idea. Suppose the Controller has to run two MapReduce jobs to expand level $i$ and $i+1$ in the tree. According to our discussion, until now it would schedule Job-1 first and then Job-2 (upper part of Figure). However, to eliminate the latency due to Setup-2, the Controller sets up Job-2 while Job-1 is still running (lower part of Figure).

To implement forward scheduling, the Controller runs a background thread which continuously keeps setting up one or more MapReduce jobs on the cluster. Once the jobs are set up, the mappers for the job wait on the Controller to send them a model file and the set of nodes to expand. When the Controller finds work on MRQ or InMemQ, it sends the work information out to the waiting mappers for a job using an RPC. With forward scheduling, lines 6 and 10 of Algorithm 8 now make RPCs rather than spawning off new threads, and the previous lines try to reserve one of the spawned MapReduces.

In practice, the Controller can forward schedule multiple jobs at the same time depending on the number of MapReduce jobs it expects to be running in parallel. A possible downside of forward scheduling is that the forward scheduling of too many jobs can result in wasted resources, where machines are waiting to receive task specifications, or in some cases receive no tasks since tree induction may be complete. Depending on availability in the cluster and the expected tree depth and ensemble type, we tune the amount of forward scheduling in the Controller.

## 6.2 Fingerprinting

Another significant source of latency that we observed in our MapReduce jobs was the cost of traversing the model: an operation performed on every mapper to determine if the training record being processed is part of the input to any node being expanded in the job. After careful examination and profiling, we found that predicate evaluations at nodes that split on unordered attributes were a bottleneck because a single predicate evaluation required multiple string comparisons, and some of our attributes were long strings, e.g., URLs. To get around this, for a predicate of the form $X \in \{v_1, v_2, \ldots v_k\}$, we fingerprint the $v_i$'s and store a hash set at the node. This simple optimization provided about 40% improvement in tree traversal costs.

## 6.3 Reliability

Deploying PLANET on a cluster of commodity machines presents a number of challenges not normally posed when running an application on a single machine. Because our clusters are shared resources, job failures due to preemption by other users is not uncommon. Similarly, job failures because of hardware issues occur occasionally. Because of the frequency of job failures, we require PLANET to have a mechanism for recovering from failures. Fortunately, the MapReduce framework provides us guarantees in terms of job completion. Therefore, we can reason about the system by considering the expansion of a set of nodes as an atomic operation and when a single MapReduce fails the Controller will simply restart the MapReduce again.

To handle the failure of the Controller, we annotate the model file with metadata marking the completion of each splitting task. Then, when the Controller fails, we start a new Controller that reads in the annotated model file generated during the failed run. Given the annotated model file, it is simple for the Controller to reconstruct the state of MRQ and InMemQ prior to any jobs which were running when the Controller failed. With MRQ, InMemQ and M, the Controller can then continue with tree induction.

Monitoring turned out to be another issue in deploying PLANET. As developers and users of the system, we often needed to be able to monitor the progress of model construction in real time. To support such monitoring, we added a dashboard to PLANET to track its currently running tasks as well as the pending tasks in MRQ and InMemQ. The dashboard collects training and validation error statistics and renders a plot of the error of the model as it grows (and offers a precision-recall curve when training a model for classification).

## 7. EXPERIMENTS

In this section we demonstrate the performance of PLA-NET on a real world learning task in computational advertising. In particular, we study the scalability of the system and the benefits obtained from the different extensions and optimizations proposed in the paper.

## 7.1 Setup

We measure the performance of PLANET on the *bounce rate prediction problem* [22, 23]. A click on an sponsored search advertisement is called a *bounce* if the click is immediately followed by the user returning to the search engine. Ads with high bounce rates are indicative of poor user experience and provide a strong signal of advertisement quality.

The training dataset (ADCORPUS) for predicting bounce rates is derived from all clicks on search ads from the Google
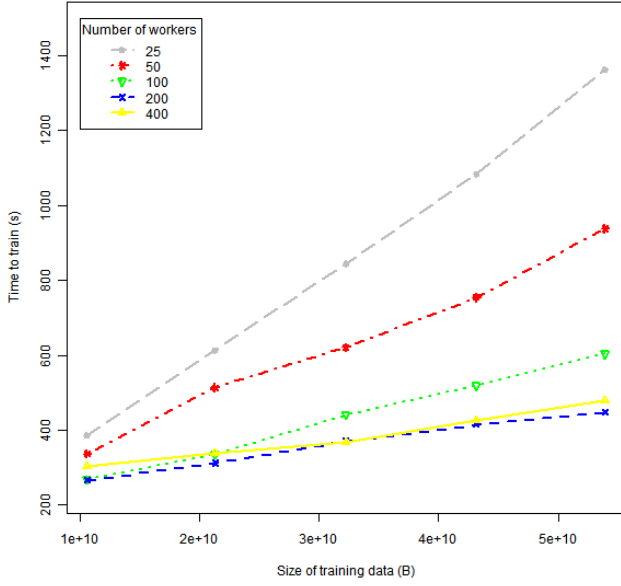
**Figure 3: Running time vs data size for various numbers of machines.**



**Figure 4: Running time vs tree depth. Note: the Sampled R curve was trained on 1/30 of the data used for the other curves.**

search engine in a particular time period. Each record represents a click labeled with whether it was bounce. A wide variety of features are considered for each click. These include the search query for the click, advertiser chosen keyword, advertisement text, estimated clickthrough rate of the ad clicked, a numeric similarity score between the ad and the landing page, and whether the advertiser keyword precisely matched the query. To improve generalization, we generalized the query and advertiser keywords into one of approximately 500 clusters, and used cluster properties as additional features. Overall, the dataset consisted of 6 categorical features varying in cardinality from 2 to 500, 4 numeric features, and 314 million records.

All of our experiments were performed on a MapReduce equipped cluster where each machine was configured to use 768MB of RAM and 1GB of hard drive space (peak utilization was < 200MB RAM and 50MB disk). Unless otherwise noted, each MapReduce job used 200 machines. A single MapReduce was never assigned more than 4 nodes for splitting and at any time a maximum of 3 MapReduce jobs were scheduled on the cluster. Running time was measured as the total time between the cluster receiving a request to run PLANET and PLANET exiting with the learned model as output. In each experiment, the first run was ignored because of the additional one-time latency to stage PLANET on the cluster. To mitigate the effects of varying cluster conditions, all the running times have been averaged over multiple runs.

To put the timing numbers that follow into perspective, we also recorded the time taken to train tree models in R using the GBM package [28]. This package requires the entire training data in memory, and hence we train on a sample of 10 million records (about 2 GB). On a machine with 8GB RAM and sufficient disk, we trained 10 trees, each at depth
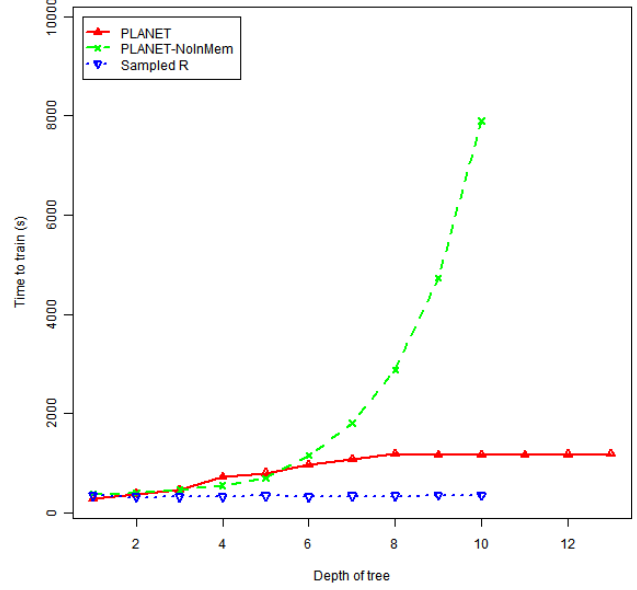
between 1 and 10. Peak RAM utilization was 6GB (average was close to 5GB). The runtime for producing the different trees varied between 315 and 358 seconds (Figure 4).

## 7.2 Results

**Scalability:** Our first experiment measures the scalability of the PLANET framework. For this experiment, we randomly split the ADCORPUS into 5 roughly equal-sized groups and trained a single depth-3 classification tree, first on a single group, then two groups and so on up to five groups. For each of these increasingly larger training datasets, we examined the effects of using between 50 and 600 machines. In this experiment, the Controller never scheduled more than 2 MapReduce jobs at a time, and was configured to schedule MR_ExpandNodes jobs only. In other words, we disabled the optimization to construct trees entirely in memory and limited forward scheduling to 1 MapReduce in order to evaluate the performance of the algorithm in a constrained (e.g. shared cluster) environment.

Figure 3 shows the results of this experiment. As expected, training time increases in proportion to the amount of training data. Similarly, adding more machines significantly decreases training time (ignoring the 400 machine curve for the moment). The most interesting observation in Figure 3 is the notion of marginal returns. When the dataset is large, adding more machine reduces costs proportionally, up to a point. For example, in our experiment, increasing the number of machines from 200 to 400 per MapReduce did not improve training time. Similarly, as the training set size decreases, the benefits of adding more machines also diminishes. In both these cases, after a certain point the overhead of adding new machines (networking overhead to watch the worker for failure, to schedule backup workers, to distribute data to the worker, and to collect results from the

worker) dominate the benefits from each machine processing a smaller chunk of data. Empirically, it appears that for our dataset the optimal number of workers is under 400.

**Benefits of MR_InMemory:** Our next experiment highlights the benefits from in memory tree completion. Here, the Controller was configured to invoke MR_InMemory for nodes whose inputs contained 10M or fewer records. The reducers in MR_InMemory used the GBM package for tree construction and were configured with 8GB RAM in order to meet the memory requirements of the package. PLANET was used to train a single classification tree of varying depths on the entire ADCORPUS.

Figure 4 shows the results. PLANET-NoInMem plots the training time when MR_InMemory is not used by the Controller. In this case training time keeps increasing with tree depth as the number of MR_ExpandNodes jobs keeps increasing. Note that even though we expand trees breadth first, the increase in training time is not linear in the depth. This happens because each MR_ExpandNodes job is configured (based on memory constraints in the mappers) to expand four nodes only. At lower levels of the tree a single MapReduce can no longer expand all nodes in a level, and hence we see a superlinear increase in training time. On the other hand, PLANET using a mix of MR_ExpandNodes and MR_InMemory scales well and training time does not increase as significantly with tree depth.

As a reference point for the PLANET running times, we also provide the running time of Sampled-R in Figure 4, which shows the running time of the GBM in-memory algorithm on a 2GB sample of ADCORPUS.
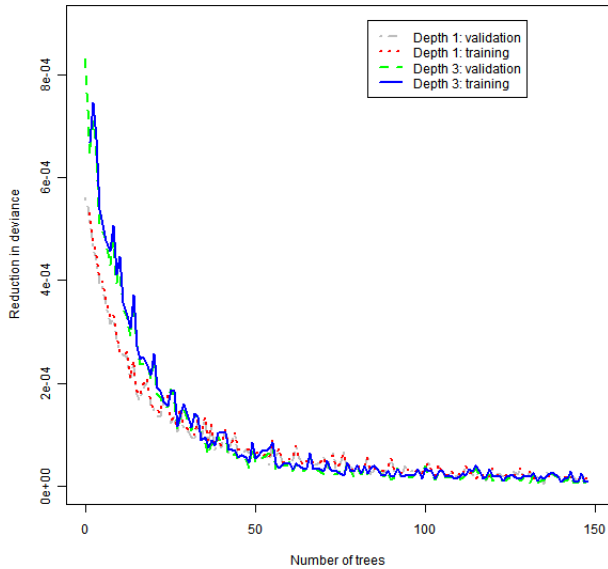


**Figure 5: Error reduction as the number of trees increases.**

**Effect of Ensembles:** The last experiment we report shows how error rates decrease in the bounce rate problem. Figure 5 shows the reduction in training and validation errors

on a 90-10 split of the ADCORPUS. The figure plots the reduction in variance as more trees are added to a boosted tree model. Two scenarios are shown – one in which the weak learners are depth one trees, and the other where the trees have depth three. For the depth-3 tree ensemble, the reduction in error is initially higher than with the depth-1 tree ensemble as expected; but, the reduction asymptotes after about 100 trees for this dataset. The PLANET dashboard updates and displays such error graphs in real time. This enables users to manually intervene and stop model training when the error converges or overfitting begins.

## 8. RELATED WORK

Scaling up tree learning algorithms to large datasets is an area of active research interest. There have been two main research directions taken by previous work: (1) centralized algorithms for large datasets on disk to avoid in-memory limitations, and (2) parallel algorithms on specific parallel computing architectures. In applying the MapReduce framework to large scale tree learning, PLANET borrows and builds upon several ideas from these previous approaches.

**Centralized Algorithms:** Notable centralized algorithms for scaling decision tree learning to large datasets include SLIQ [26], CLOUDS [1], RAINFOREST [17], and BOAT [16]. SLIQ uses strategies like pre-sorting and attribute lists in breadth-first tree-growing to enable learning from large training data on disk. While PLANET does not use pre-sorting or attribute lists, it grows the tree breadth-first like SLIQ. The key insight in RAINFOREST is that the splitting decision at a tree node needs a compact data structure of sufficient statistics (called AVC group in the paper), which in most cases can be fit in-memory. PLANET similarly maintains sufficient statistics on mappers during MR_ExpandNodes. CLOUDS samples the split points for numeric attributes and uses an estimation step to find the best split point, resulting in lower computation and I/O cost compared to other tree learning algorithms like C4.5. For efficient estimation of the best split, PLANET uses equidepth histograms of ordered attributes to estimate split points. Finally, BOAT uses statistical sampling to construct a tree based on a small subset of the whole data and then does corrections to the tree based on estimated differences compared to the actual tree learned on the whole data. In comparison, PLANET builds the tree from the whole data directly.

**Parallel Algorithms:** Numerous approaches for parallelizing tree learning have been proposed. [27] contains an excellent survey of existing approaches, along with the motivations for large scale tree learning. Bradford et al. [3] discuss how the C4.5 decision tree induction algorithm can be effectively parallelized in the ccNUMA parallel computing platform. It also mentions other parallel implementations of decision trees, namely SLIQ, SPRINT, and ScalParC for message-passing systems, and SUBTREE, MWK and MLC++ for SMPs. Most of these algorithms have been developed for specific parallel computing architectures, many of which have specific advantages, e.g., shared memory to avoid replicating or communicating the whole dataset among the processors. In comparison, PLANET is based on the MapReduce platform that uses commodity hardware for massive-scale parallel computing.

For deciding the split points of attributes, SPRINT [31] uses attribute lists like SLIQ. Each processor is given a sub-list of each attribute list, corresponding to the instance indices in the data chunk sent to the processor. While computing good split points, each processor determines the gains over the instances assigned to that processor for each ordered attribute, and sends the master a portion of the statistics needed to determine the best split. However, this requires an all-to-all broadcast of instance ids at the end. PLANET takes a simpler and more scalable approach – instead of considering all possible split points, it computes a representative subset of the splits using approximate histograms, after which the selection of the best split can be done using only one MapReduce job (details in Section 4.1).

ScalParC [21], which builds on SLIQ and SPRINT, also splits each attribute list into multiple parts and assigns each part to a processor. However, rather than building the tree in a depth-first manner (as done by C4.5, MLC++, etc.), it does a breadth-first tree growth like SLIQ (and PLANET) to prevent possible load imbalance in a parallel computing framework.

Other notable techniques for parallel tree learning include: (1) parallel decision tree learning on a SMP architecture based on attribute scheduling among processors, including task pipelining and dynamic load balancing for speedup [33]; (2) meta-learning schemes that train multiple trees in parallel along with a final arbiter tree that combines their predictions [10]; (3) distributed learning of trees by boosting, which operates over partitions of a large dataset that are exchanged among the processors [24]; (4) the SPIES algorithm, which combines the AVC-group idea of RAINFOR-EST with effective sampling of the training data to obtain a communication- and memory-efficient parallel tree learning method [19]; (5) a distributed tree learning algorithm that uses only 20% of the communication cost to centralize the data, but achieves 80% of the accuracy of the centralized version [18].

On the theoretical side, Caragea et al. [7] formulated the problem of learning from distributed data and showed different algorithm settings for learning trees from distributed data, each of which is provably exact, i.e., they give the same results as a tree learned using all the data in a centralized setting. Approximate algorithms for parallel learning of trees on streaming data have also been recently proposed [2, 20].

**MapReduce in Machine Learning:** In recent years, some learning algorithms have been implemented using the MapReduce framework. Chu et al. [11] give an excellent overview of how different popular learning algorithms (e.g., locally weighted linear regression, naïve Bayes classification, Gaussian discriminant analysis, k-means, logistic regression, neural networks, principal component analysis, independent component analysis, expectation maximization, support vector machines) can be effectively solved in the MapReduce framework. However, these algorithms have all been implemented using a shared-memory multi-processor architecture. Our focus is on scaling learning algorithms (especially ensemble tree learning) to massive datasets using a MapReduce framework deployed on commodity hardware.

## 9. CONCLUSIONS

We have presented PLANET, a framework for large-scale tree learning using a MapReduce cluster. We are currently applying PLANET to problems within the sponsored search domain. Our experience is that the system scales well and performs reliably in this context, and we expect results would be similar in a variety of other domains involving large scale learning problems. Our initial goal in building PLANET was to develop a scalable tree learner with accuracy comparable to a traditional in-memory algorithm, but capable of handling much more training data. We believe our experience in building and deploying PLANET provides lessons in using MapReduce for other non-trivial mining and data processing tasks. The strategies we developed for handling tree learning should be applicable to other problems requiring multiple iterations, each requiring one or more applications of MapReduce.

For future work, our short term focus is to extend the functionality of PLANET in various ways to support more learning problems at Google. For example, we intend to support split metrics other than those based on variance. We also intend to investigate how intelligent sampling schemes might be used in conjunction with the scalability offered by PLANET. Other future plans include extending the implementation to handle multi-class classification and incremental learning.

## 10. REFERENCES

[1] K. Alsabti, S. Ranka, and V. Singh. Clouds: A decision tree classier for large datasets. Technical report, University of Florida, 1998.

[2] Y. Ben-Haim and E. Yom-Tov. A streaming parallel decision tree algorithm. In *Large Scale Learning Challenge Workshop at the International Conference on Machine Learning (ICML)*, 2008.

[3] J. P. Bradford, J. A. B. Fortes, and J. Bradford. Characterization and parallelization of decision tree induction. Technical report, Purdue University, 1999.

[4] L. Breiman. Bagging predictors. *Machine Learning Journal*, 24(2):123–140, 1996.

[5] L. Breiman. Random forests. *Machine Learning Journal*, 45(1):5–32, 2001.

[6] L. Breiman, J. H. Friedman, R. Olshen, and C. Stone. *Classification and Regression Trees*. Wadsworth and Brooks, 1984.

[7] D. Caragea, A. Silvescu, and V. Honavar. A framework for learning from distributed data using sufficient statistics and its application to learning decision trees. *International Journal of Hybrid Intelligent Systems*, 1(1–2):80–89, 2004.

[8] R. Caruana, N. Karampatziakis, and A. Yessenalina. An empirical evaluation of supervised learning in high dimensions. In *International Conference on Machine Learning (ICML)*, pages 96–103, 2008.

[9] R. Caruana and A. Niculescu-Mizil. An empirical comparison of supervised learning algorithms. In

*International Conference on Machine Learning (ICML)*, pages 161–168, 2006.

[10] P. K. Chan and S. J. Stolfo. Toward parallel and distributed learning by meta-learning. In *Workshop on Knowledge Discovery in Databases at the Conference of Association for the Advancement of Artificial Intelligence (AAAI)*, pages 227–240, 1993.

[11] C.-T. Chu, S. K. Kim, Y.-A. Lin, Y. Yu, G. Bradski, A. Y. Ng, and K. Olukotun. Map-reduce for machine learning on multicore. In *Advances in Neural Information Processing Systems (NIPS) 19*, pages 281–288, 2007.

[12] J. Dean and S. Ghemawat. Mapreduce: Simplified data processing on large clusters. In *Symposium on Operating System Design and Implementation (OSDI)*, 2004.

[13] R. O. Duda, P. E. Hart, and D. G. Stork. *Pattern Classification*. Wiley, New York, second edition, 2001.

[14] Y. Freund and R. E. Schapire. Experiments with a new boosting algorithm. In *International Conference on Machine Learning (ICML)*, pages 148–156, 1996.

[15] J. H. Friedman. Greedy function approximation: A gradient boosting machine. *Annals of Statistics*, 29(5):1189–1232, 2001.

[16] J. Gehrke, V. Ganti, R. Ramakrishnan, and W.-Y. Loh. BOAT – Optimistic decision tree construction. In *International Conference on ACM Special Interest Group on Management of Data (SIGMOD)*, pages 169–180, 1999.

[17] J. Gehrke, R. Ramakrishnan, and V. Ganti. Rainforest - A framework for fast decision tree construction of large datasets. In *International Conference on Very Large Data Bases (VLDB)*, pages 416–427, 1998.

[18] C. Giannella, K. Liu, T. Olsen, and H. Kargupta. Communication efficient construction of decision trees over heterogeneously distributed data. In *International Conference on Data Mining (ICDM)*, pages 67–74, 2004.

[19] R. Jin and G. Agrawal. Communication and memory efficient parallel decision tree construction. In *SIAM Conference on Data Mining (SDM)*, pages 119–129, 2003.

[20] R. Jin and G. Agrawal. Efficient decision tree construction on streaming data. In *SIGKDD Conference on Knowledge Discovery and Data Mining (KDD)*, pages 571–576, 2003.

[21] M. Joshi, G. Karypis, and V. Kumar. Scalparc: A new scalable and efficient parallel classification algorithm for mining large datasets. In *International Parallel Processing Symposium (IPPS)*, pages 573–579, 1998.

[22] A. Kaushik. Bounce rate as sexiest web metric ever. MarketingProfs, August 2007. http://www.marketingprofs.com/7/bounce-rate-sexiest-web-metric-ever-kaushik.asp?sp=1.

[23] A. Kaushik. Excellent analytics tip 11: Measure effectiveness of your web pages. Occam's Razor (blog), May 2007. http://www.kaushik.net/avinash/2007/05/excellent-analytics-tip-11-measure-effectiveness-of-your-web-pages.html.

[24] A. Lazarevic. The distributed boosting algorithm. In *SIGKDD Conference on Knowledge Discovery and Data Mining (KDD)*, pages 311–316, 2001.

[25] G. S. Manku, S. Rajagopalan, and B. G. Lindsay. Random sampling techniques for space efficient online computation of order statistics of large datasets. In *International Conference on ACM Special Interest Group on Management of Data (SIGMOD)*, pages 251–262, 1999.

[26] M. Mehta, R. Agrawal, and J. Rissanen. Sliq: A fast scalable classifier for data mining. In *International Conference on Extending Data Base Technology (EDBT)*, pages 18–32, 1996.

[27] F. Provost and U. Fayyad. A survey of methods for scaling up inductive algorithms. *Data Mining and Knowledge Discovery*, 3:131–169, 1999.

[28] G. Ridgeway. Generalized boosted models: A guide to the gbm package. http://cran.r-project.org/web/packages/gbm, 2006.

[29] L. Rokach and O. Maimon. *Data Mining with Decision Trees: Theory and Applications*. World Scientific Publishing Company, 2008.

[30] D. Sculley, R. Malkin, S. Basu, and R. J. Bayardo. Predicting bounce rates in sponsored search advertisements. In *SIGKDD Conference on Knowledge Discovery and Data Mining (KDD)*, pages 1325–1334, 2009.

[31] J. C. Shafer, R. Agrawal, and M. Mehta. Sprint: A scalable parallel classifier for data mining. In *International Conference on Very Large Data Bases (VLDB)*, pages 544–555, 1996.

[32] V. N. Vapnik. *The Nature of Statistical Learning Theory*. Springer-Verlag, Berlin, 1995.

[33] M. J. Zaki, C.-T. Ho, and R. Agrawal. Parallel classification for data mining on shared-memory multiprocessors. In *International Conference on Data Engineering (ICDE)*, pages 198–205, 1999.