

See discussions, stats, and author profiles for this publication at: <https://www.researchgate.net/publication/372479561>

A Scalable Tree Boosting System: XG Boost

Article · January 2020

DOI: 10.22259/2349-476X.0712005

CITATIONS

50

READS

3,105

3 authors, including:



Mounika Nalluri

Cigna

29 PUBLICATIONS 167 CITATIONS

[SEE PROFILE](#)



Nageswara rao Eluri

R.V.R. & J.C. College of Engineering

25 PUBLICATIONS 116 CITATIONS

[SEE PROFILE](#)

A Scalable Tree Boosting System: XG Boost

Mounika Nalluri¹, Mounika Pentela¹, Nageswara Rao Eluri²

¹CSE Department, Malineni Lakshmaiah Women's, Engineering College, Guntur, AP

²Associate. Professor, CSE Department, Malineni Lakshmaiah Women's, Engineering College, Guntur, AP

ABSTRACT

Tree boosting is a well-known and effective machine learning method. We offer XGBoost, an end-to-end scalable tree boosting system that is used in this paper. It is heavily used by data scientists to achieve cutting-edge results. addressing a wide range of machine learning issues We propose a fresh approach. A sparsity-aware approach, as well as a weighted quantile sketch for approximation tree learning, are employed for sparse data. However, we give insights on cache utilisation patterns, data compression, and sharding to construct a scalable tree boosting system. By combining these concepts, XG Boost can scale beyond billions. a number of cases where the resources used are much less than those used by present systems

Keywords: Large-Scale Machine Learning

EXAMPLE OF INTRODUCTION

Machine learning and data-driven approaches are becoming increasingly important in a variety of industries. Smart spam classifiers learn to safeguard our email by analyzing vast quantities of spam data and user feedback; advertising systems learn to connect the appropriate adverts with the right people.

In addition to being utilized as a standalone predictor, it is integrated into real-world production workflows for ad click through rate prediction [15]. Last but not least, it is the de facto ensemble approach.

This study introduces XGBoost, a scalable machine learning framework for tree boosting. The system is available for download and is open source². The system's significance has been widely recognised in a number of machine learning and data mining problems. Consider the challenges offered by Kaggle, a machine learning competition website. In 2015, XGBoost was used in 17 of the 29 winning challenge solutions published on Kaggle's blog. Only XGBoost was used to train the model in eight of these solutions, whereas the remainder used a combination of XGBoost and neural nets in ensembles. In 11 of the solutions, deep neural nets, the second most prevalent method, were applied. Every winning team in the top ten in the KDDCup 2015 used XGBoost, demonstrating the system's effectiveness. Ensemble techniques,

according to the winning teams, only marginally outperform a well-configured XGBoost [1].

These findings demonstrate that our strategy is effective. For a wide range of situations, cutting-edge results are available. One of the issues addressed in these winning solutions is store sales. Online text consumer behaviour prediction; motion detection; ad click through rate prediction; malware classification; Product classification; hazard risk forecasting; and huge online course dropout rate forecasting Domain-dependent data analysis and feature engineering are important, but they're not the only ones.

The fact that XGBoost is the learner's unanimous choice indicates how important these solutions are. tree augmentation and our system. The system's scalability in all scenarios is the most important aspect of XGBoost's performance. More than one user can be accommodated by the system. It is ten times faster than existing popular solutions on a single system. Machine learning may scale to billions of examples in distributed or parallel computing. Configurations with a memory limit. The scalability of XGBoost is due to a variety of factors.

As a result of a number of significant system and algorithmic enhancements On a single machine, data scientists can process hundreds of millions of examples. Finally, integrating

these ideas to create an end-to-end system that scales to even more data while using the smallest amount of cluster resources is even more intriguing. The following are the paper's key contributions. We offer a theoretically justified weighted quantile sketch for efficient proposal calculation, and we construct and build an end-to-end tree boosting system that is highly scalable.

A one-of-a-kind solution for real-life scenarios.

This enables academics and data scientists to develop more effective tree boosting algorithms [7, 8]. We also propose a regularised learning objective, which we will include for completeness' sake. The remainder of the paper is formatted as follows. We'll go over tree boosting in Section 2 and offer a regularised target. The split finding methods in Section 3 and the system design in Section 4 are then detailed, with experimental results supplied as needed to provide quantitative rationale for each optimization. Related work is discussed in Section 5. A new sparsity-aware parallel tree learning approach is presented. We recommend a cache-aware block structure for out-of-core tree learning. While there have been some prior work on parallel tree boosting [22, 23, 19], novel directions including out-of-core computation, cache-aware learning, and sparsity-aware learning have yet to be explored. More importantly, an end-to-end system that has all of these capabilities gives a comprehensive solution.

A TREE BOOSTING NUTSHELL

We'll look at how to improve a gradient tree in this section. The formula is based on an idea that has just been discovered in gradient boosting research. The second order approach was developed by Friedman et al. [12]. We make a couple of changes to the regularised aim that have proven to be helpful in practise.

$$\mathcal{L}(\phi) = \sum_i l(\hat{y}_i, y_i) + \sum_k \Omega(f_k) \quad (2)$$

$$\text{where } \Omega(f) = \gamma T + \frac{1}{2} \lambda \|w\|^2$$

L is a differentiable convex loss function, measures the difference between the forecast \hat{y}_i and the target y_i . The model's complexity is penalised in the second term. (Or, to put it another way, the regression tree functions.) To avoid skewed outcomes, the additional

Learning Objective (Regularized)

($|D| = n, x_i \in \mathbb{R}^m, y_i \in \mathbb{R}$) $D = (x_i, y_i)_{i=1}^n$ for a data set with n examples and m attributes. To solve the problem, a tree ensemble model (shown in Fig. 1) employs K additive functions ($m, y_i \in \mathbb{R}$) predicting the outcome.

$$\hat{y}_i = \phi(x_i) = \sum_{k=1}^K f_k(x_i), \quad f_k \in \mathcal{F}, \quad (1)$$

where $F = \{f(x) = wq(x) \mid q \in \mathcal{Q}\}$ and $\mathcal{Q} = \{q(x) = wq(x) \mid q \in \mathcal{Q}\}$ (also known as CART). q represents the structure of any tree that maps an example to q , the leaf index that corresponds to it. T tree refers to the number of leaves on the tree. Each f_k denotes a unique tree structure. w and q leaf weights. Each regression model, unlike decision trees, is unique. A tree with a continuous score on each leaf is used.

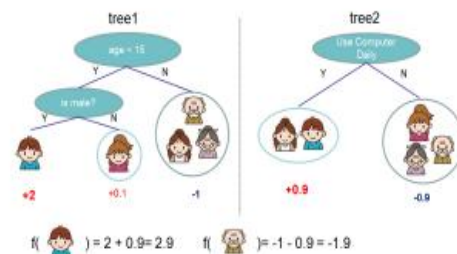


Figure 1: Tree Ensemble Model. The final prediction for a given example is the sum of predictions from each tree.

Put it in the leaves, and then add up all the scores in the appropriate leaves to get the final prediction (given by w). To learn the set of functions used in the model, we minimise the following regularised goal. The score on the i -th leaf is represented by w_i . Let's look at an illustration. will use the decision rules in the trees to categorise (provided by q).

regularisation term aids in the smoothing of the final learned weights. over-fitting. Intuitively, the regularised goal will lead to the selection of a model that employs simple and predictive functions. A similar regularisation method is used in the Regularized Greedy

Forest (RGF) [25] model. Our mission and goal The learning algorithm linked with RGF is less complicated. Parallelization is also easier. The goal reverts to its previous state when the regularisation option is set to zero.

Gradient Tree Enhancement

Gradient Tree Boosting is a technique for increasing the size of a tree. Because the tree ensemble model in Eq. (2) uses functions as parameters, it can't be optimised in Euclidean space with traditional optimization methods. Rather, the paradigm is presented in a way that is both simple and effective. There are two types of subtractive and additive functions. Let $y(t)$ be the one in a formal sense. We will forecast the i -th case at the t -th iteration. You'll need to add f_t to the following aim to minimise it.

$$\mathcal{L}^{(t)} = \sum_{i=1}^n l(y_i, \hat{y}_i^{(t-1)} + f_t(\mathbf{x}_i)) + \Omega(f_t)$$

This implies we add the foot that Eq thinks will improve our model the most (2). Second-order approximation can be used to swiftly optimise the target in general. The twelfth setting This implies we add the foot that Eq thinks will improve our model the most (2). Second-order approximation can be used to quickly optimise the goal in the wide case [12].

Define $I_j = \{i | q(\mathbf{x}_i) = j\}$ as the instance set of leaf j . We can rewrite Eq (3) by expanding Ω as follows

$$\begin{aligned} \tilde{\mathcal{L}}^{(t)} &= \sum_{i=1}^n [g_i f_t(\mathbf{x}_i) + \frac{1}{2} h_i f_t^2(\mathbf{x}_i)] + \gamma T + \frac{1}{2} \lambda \sum_{j=1}^T w_j^2 \\ &= \sum_{j=1}^T [(\sum_{i \in I_j} g_i) w_j + \frac{1}{2} (\sum_{i \in I_j} h_i + \lambda) w_j^2] + \gamma T \end{aligned} \quad (4)$$

For a fixed structure $q(\mathbf{x})$, we can compute the optimal weight w_j^* of leaf j by

$$w_j^* = - \frac{\sum_{i \in I_j} g_i}{\sum_{i \in I_j} h_i + \lambda}, \quad (5)$$

and calculate the corresponding optimal value by

$$\tilde{\mathcal{L}}^{(t)}(q) = - \frac{1}{2} \sum_{j=1}^T \frac{(\sum_{i \in I_j} g_i)^2}{\sum_{i \in I_j} h_i + \lambda} + \gamma T. \quad (6)$$

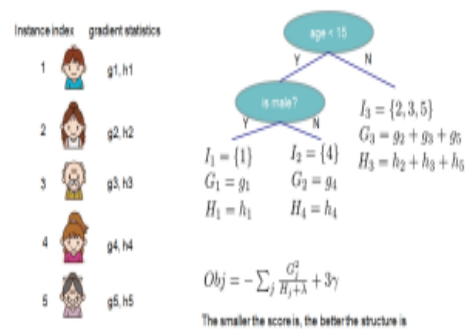
the left and right instance sets are IL and IR , respectively. There are left and right nodes after the break. If we set I equal to IL IR , we get The loss reduction is given by after the split.

$$\mathcal{L}^{(t)} \simeq \sum_{i=1}^n [l(y_i, \hat{y}_i^{(t-1)}) + g_i f_t(\mathbf{x}_i) + \frac{1}{2} h_i f_t^2(\mathbf{x}_i)] + \Omega(f_t)$$

where $g_i = \partial_{\hat{y}_i^{(t-1)}} l(y_i, \hat{y}_i^{(t-1)})$ and $h_i = \partial_{\hat{y}_i^{(t-1)}}^2 l(y_i, \hat{y}_i^{(t-1)})$

On the loss function, there exist first and second order gradient statistics. At step t , we can eliminate the constant terms to obtain the following reduced objective.

$$\tilde{\mathcal{L}}^{(t)} = \sum_{i=1}^n [g_i f_t(\mathbf{x}_i) + \frac{1}{2} h_i f_t^2(\mathbf{x}_i)] + \Omega(f_t) \quad (3)$$



Calculation of the Structure Score (Figure 2). To calculate the quality score, simply add the gradient and second order gradient statistics on each leaf, then use the scoring algorithm.

The quality of a tree structure q can be determined using Eq (6) as a criterion. This score is calculated in the same way as the impurity score for analysing decision trees. for a more varied range of goal functions Figure 2 shows What formula is used to calculate this score?

Normally, listing all of the possible options q is difficult. Starting with a tree structure, a greedy method is used. Instead of starting with a single leaf, iteratively adds branches to the tree. Assume that

$$\mathcal{L}_{split} = \frac{1}{2} \left[\frac{(\sum_{i \in IL} g_i)^2}{\sum_{i \in IL} h_i + \lambda} + \frac{(\sum_{i \in IR} g_i)^2}{\sum_{i \in IR} h_i + \lambda} - \frac{(\sum_{i \in I} g_i)^2}{\sum_{i \in I} h_i + \lambda} \right] - \gamma \quad (7)$$

In practice, this formula is commonly used to evaluate split candidates.

Shrinkage and Column Subsampling

In addition to the regularised objective mentioned in Sec. 2.1, two additional methods are used to prevent overfitting. The first strategy, shrinkage, was proposed by Friedman [11]. The size of freshly added weights shrinks by a factor of two due to shrinkage. following each step of the tree boosting process It's comparable to the concept of a learning rate. The influence of shrinking towards chiastic optimization is reduced. Each tree is unique, and room is left in the model for future trees to improve it. Subsampling by column is the second technique (feature). This method is used in Random Forest [4,].

Algorithm 1: Exact Greedy Algorithm for Split Finding

Input: I , instance set of current node

Input: d , feature dimension

 $gain \leftarrow 0$
$$G \leftarrow \sum_{i \in I} g_i, H \leftarrow \sum_{i \in I} h_i$$
for $k = 1$ to m do
$$G_L \leftarrow 0, H_L \leftarrow 0$$
for j in sorted(I , by \mathbf{x}_{jk}) do
$$| \quad G_L \leftarrow G_L + g_j, \quad H_L \leftarrow H_L + h_j$$
$$G_R \leftarrow G - G_L, \quad H_R \leftarrow H - H_L$$
$$score \leftarrow \max(score, \frac{G_L^2}{H_L + \lambda} + \frac{G_R^2}{H_R + \lambda} - \frac{G^2}{H + \lambda})$$

end

end

Output: Split with max score

Algorithm 2: Approximate Algorithm for Split Finding

for $k = 1$ to m do

Propose $S_k = \{s_{k1}, s_{k2}, \dots, s_{kl}\}$ by percentiles on feature k .

Proposal can be done per tree (global), or per split(local).

end

for $k = 1$ to m do
$$G_{kv} \leftarrow \sum_{i \in I_j | e_i \sim v, \dots, e_i \sim v} g_i$$
$$H_{kv} \leftarrow \sum_{j \in \{j | s_{k,v} > \mathbf{x}_{jk} > s_{k,v-1}\}} h_j$$

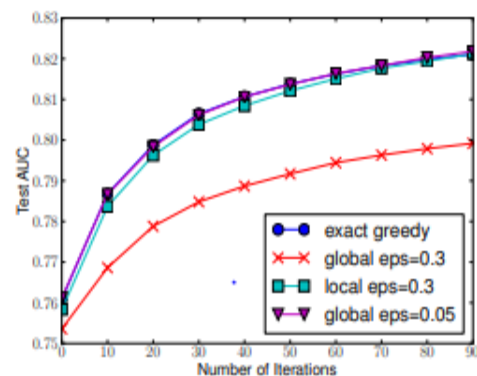
end

Follow same step as in previous section to find max score only among proposed splits.

According to user input, column sub-sampling lowers over-fitting even more than row sub-sampling. Row sub sampling is a sample technique that has been around for quite some time (which is also supported). The parallel approach, which will be explained later, is also sped up by using column sub-samples.

ALGORITHMS FOR SPLIT FINDING

Exact Greedy Algorithm (Basic)

[illegible]

[13], For gradient boosting, it's featured in the commercial product Tree Net 4, however it's not yet included in any open source packages.

Figure 3 illustrates a comparison of test AUC convergence on the Higgs 10M dataset. The eps parameter represents the precision of the approximation drawing. This effectively translates to $1 / \text{eps}$ buckets in the proposal. Local proposals require fewer bins, as we've discovered. It refines split candidates due to the fact that it refines split candidates.

A tribute to the past To provide efficient gradient tree boosting in both of these scenarios, an approximation technique is necessary. We describe a preliminary framework in this paper that parallels ideas proposed in earlier literatures [17, 2, 22].

Algorithm 2 To summarise, the algorithm first suggests candidate splitting locations based on percentiles of feature distribution (a specific criteria will be given in Sec. 3.3). The programme next separates the continuous features into buckets based on these candidate points and collects the information. It selects the best selection from a list of possibilities.

Depending on when the suggestion is made, the algorithm has two versions. The global variant provides all possible splits during the first phase of tree construction and uses the same split finding recommendations at all levels. The native variety re-produces after each split. The entire globe. There are fewer proposal phases in this method than in the local technique. In most cases, though, more candidate points are required. This is a worldwide proposition split because candidates are not refined after each round. The local proposal refines the competitors after splits. It may be better suited to trees with a more extensive root system. Multiple approaches were compared on a Higgs boson dataset. The situation is depicted in Figure 3. We've learned that the local suggestion is unquestionably viable. There are fewer candidates needed. The following is a general proposal: It will be as accurate as the poll if there are enough candidates.

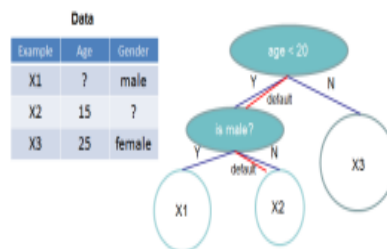


Figure4. Default tree structure and directions. When a feature required for the split is lacking, an example will be classified in the default direction.

ta. Formally, let multi-set $\mathcal{D}_k = \{(x_{1k}, h_1), (x_{2k}, h_2) \dots (x_{nk}, h_n)\}$ represent the k -th feature values and second order gradient statistics of each training instances. We can define a rank functions $r_k : \mathbb{R} \rightarrow [0, +\infty)$ as

$$r_k(z) = \frac{1}{\sum_{(x,h) \in \mathcal{D}_k} h} \sum_{(x,h) \in \mathcal{D}_k, x < z} h, \quad (8)$$

which represents the proportion of instances whose feature value k is smaller than z . The goal is to find candidate split points $\{s_{k1}, s_{k2}, \dots, s_{kl}\}$, such that

$$|r_k(s_{k,j}) - r_k(s_{k,j+1})| < \epsilon, \quad s_{k1} = \min_i x_{ik}, \quad s_{kl} = \max_i x_{ik}. \quad (9)$$

Here ϵ is an approximation factor. Intuitively, this means that there is roughly $1/\epsilon$ candidate points. Here each data point is weighted by h_i . To see why h_i represents the weight, we can rewrite Eq (3) as

$$\sum_{i=1}^n \frac{1}{2} h_i (f_i(\mathbf{x}_i) - g_i/h_i)^2 + \Omega(f_t) + \text{constant},$$

A Scalable Tree Boosting System: XG Boost

This is exactly weighted squared loss, with labels g_i/h_i and weights h_i . In huge datasets, finding candidate splits that match the criteria is tough. When all instances have equal weights, a method known as quantile sketch [14, 24] solves the problem. But there isn't any such thing.

There is already a quantile drawing for weighted datasets. As a result, the vast majority of existing approximate algorithms resorted to sorting a random portion of data

with a chance of being helpful failure or heuristics that aren't mathematically solid. To solve this problem, we designed a revolutionary distributed system. a weighted quantile sketch algorithm that can work with weighted quantiles data with a theoretically backed guarantee that can be demonstrated. The general idea is to provide a data structure that can be merged and pruned. Each of these operations has been proved to maintain a specified level of precision and quality.

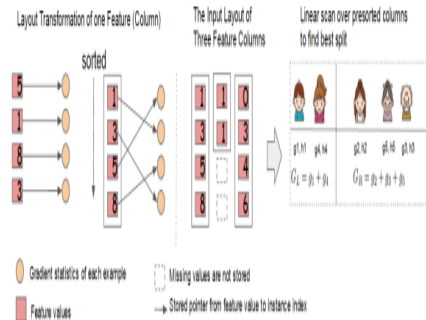


Figure 6: Block structure for parallel learning. Each column in a block is sorted by the corresponding feature value. A linear scan over one column in the block is sufficient to enumerate all the split points.

Algorithm 3: Sparsity-aware Split Finding

Input: I , instance set of current node

Input: $I_k = \{i \in I | x_{ik} \neq \text{missing}\}$

Input: d , feature dimension

Also applies to the approximate setting, only collect statistics of non-missing entries into buckets

$\text{gain} \leftarrow 0$

$G \leftarrow \sum_{i \in I} g_i, H \leftarrow \sum_{i \in I} h_i$

for $k = 1$ **to** m **do**

// enumerate missing value goto right

$G_L \leftarrow 0, H_L \leftarrow 0$

for j in $\text{sorted}(I_k, \text{ascend order by } x_{jk})$ **do**

$G_L \leftarrow G_L + g_j, H_L \leftarrow H_L + h_j$

$G_R \leftarrow G - G_L, H_R \leftarrow H - H_L$

$\text{score} \leftarrow \max(\text{score}, \frac{G_L^2}{H_L + \lambda} + \frac{G_R^2}{H_R + \lambda} - \frac{G^2}{H + \lambda})$

end

// enumerate missing value goto left

$G_R \leftarrow 0, H_R \leftarrow 0$

for j in $\text{sorted}(I_k, \text{descent order by } x_{jk})$ **do**

$G_R \leftarrow G_R + g_j, H_R \leftarrow H_R + h_j$

$G_L \leftarrow G - G_R, H_L \leftarrow H - H_R$

$\text{score} \leftarrow \max(\text{score}, \frac{G_L^2}{H_L + \lambda} + \frac{G_R^2}{H_R + \lambda} - \frac{G^2}{H + \lambda})$

end

end

Output: Split and default directions with max gain

The default direction has already been chosen. There are two possibilities for default direction in each branch. The data is utilised to figure

out which default directions are the best. The algorithm is depicted in Alg. 3. The most crucial adjustment is to just visit the locations

that aren't missing. I've made a few entries The algorithm presented here is used to deal with non-presence. as a value that is lacking and learns the best method to deal with it There are no values. In both circumstances, the same algorithm can be employed. when the absence corresponds to a user-supplied value by limiting the enumeration to consistent solutions To the best of our knowledge, most existing tree learning algorithms are either intended for dense data or require it. Specific procedures are required in some situations, such as categorical encoding. XGBoost manages all sparsity patterns in a coordinated

manner. More importantly, our approach makes use of , our method takes advantage of the sparsity to keep calculation time constant as the number of non-missing entries in the input decreases. Figure 5 shows a sparsity aware versus a naive strategy on an Allstate-10K dataset (description of dataset given in Sec. 6). The sparsity conscious method is proven to be 50 times faster than the traditional method. The naive version is the quickest. This underscores the gravity of the situation. The sparsity aware algorithm is a sort of algorithm that takes data sparsity into account.

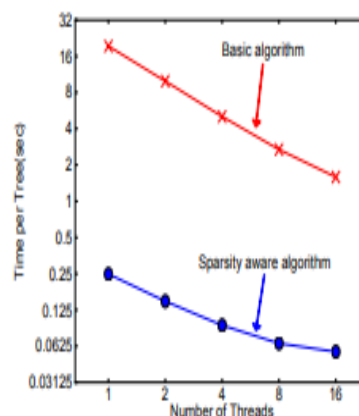


Figure 5 depicts the impact of the sparsity aware algorithm on Allstate-10K. The dataset is sparse due to one-hot encoding. The sparsity aware algorithm is a way of estimating a given's probability that is 50 times faster than the original. That does not account for sparsity.

DEVELOPMENT OF THE SYSTEM

Column Block for Group Learning

Getting to know the trees is the most time-consuming component of tree learning. Sort the information into a logical sequence. We propose sorting the data and storing it in in-memory units to save money on transit. It was given the name Block by us. The data for each block is saved in the compressed column (CSC) format, and each column is sorted in the compressed column (CSC) format. by the related trait's value This is how the info is entered. Before training, it just needs to be computed once, and it can be done several times. was reused in later versions In the precise greedy algorithm, we save the entire dataset. Run the split search procedure by linearly scanning over the pre-sorted entries in

a single block. We are the ones who have created the chasm. Because all of the leaves must be located together, the block must be scanned once. In each leaf branch, data about the split candidates will be compiled. As a result, a single scan of the block will capture information about split candidates in all leaf branches. Figure 6 shows how we transform a dataset into a graph. Format and find the optimal split using the block structure. The block structure is very useful when using approximate approaches. Many blocks can be used in this case. Each block represents a subset of the rows in the dataset. Different blocks can be spread between machines or kept on disc in an out-of-core structure. Using the information that has been sorted The discovery of the quantile . Using the information that has been sorted As a result of the structure, the quantile discovery step becomes a linear scan. over the previously sorted columns This is especially beneficial for local proposal algorithms, which create candidates at each branch on a regular basis. The binary search evolves into a linear time merging procedure in histogram aggregation.

A Scalable Tree Boosting System: XG Boost

The collecting of statistics for each column can be parallelized, allowing us to locate splits

using a parallel technique

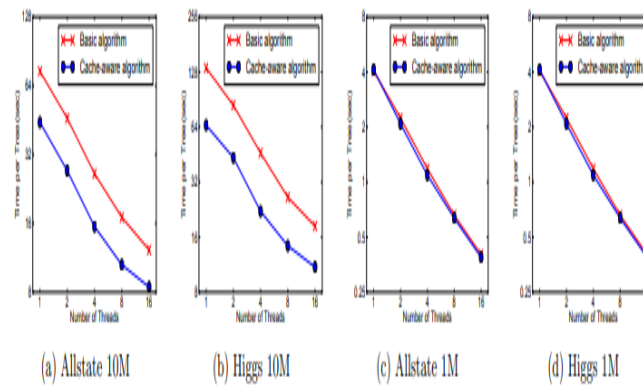


Figure 7 illustrates the effect of cache-aware pre fetching in an exact greedy algorithm. On large datasets, we discover that the cache-miss effect has an impact on performance (10

million instances). When the dataset is huge, cache aware pre fetching boosts performance by a factor of two.

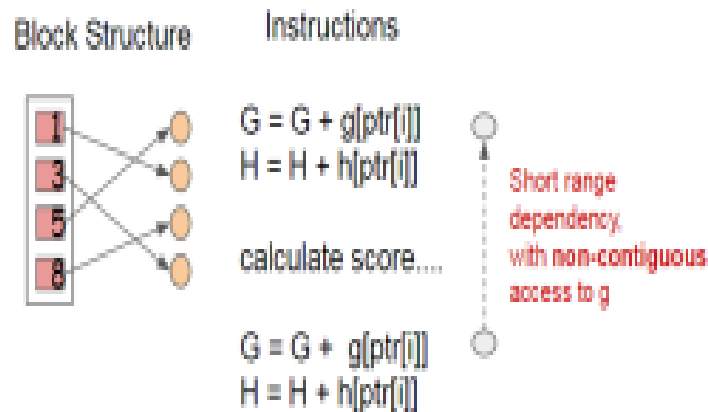


Figure8. A pattern of short-term data reliance that can create a stall due to a cache miss.

Time Complexity Analysis Let d represent the maximum depth of the tree and K represent the total number of trees. For the exact greedy method, the time complexity of the original sparse conscious approach is $O(Kdkxk_0 \log n)$. kxk_0 is the number of non-missing elements in the training data. On the other hand, tree boosting on the block structure costs merely $O(Kdkxk_0 + kxk_0 \log n)$. The $O(kxk_0 \log n)$ one-time preprocessing cost can be amortised. According to this analysis, the block structure helps save an additional $\log n$ factor, which is essential when n is large. For the approximate technique, the time complexity of the original binary search approach is $O(Kdkxk_0 \log q)$. q is the number of proposal candidates in the dataset. Even though q is usually between 32 and 100, the \log factor still adds overhead. taking advantage.

Access with Cache Awareness

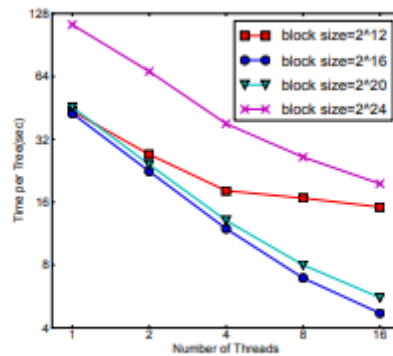
While the proposed block structure helps to reduce the complexity of split finding computations, the unique technique requires indirect fetches of gradient statistics per row index because these data are fetched in feature order. This is a non-consecutive memory access. In a naive implementation of split enumeration, the accumulation and non-continuous memory fetch operations are both read/write dependent (see Fig. 8). Split finding is slowed when the gradient statistics do not fit into the CPU cache and a cache miss occurs.

The exact greedy method can be solved using a cache-aware pre fetching approach. We allocate an internal buffer in each thread, fetch the gradient statistics into it, and then

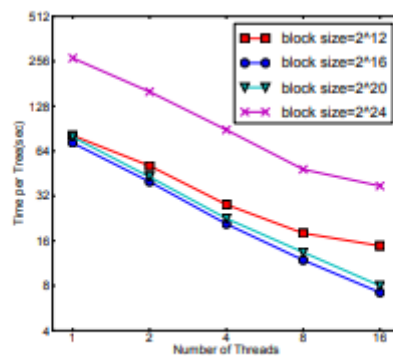
A Scalable Tree Boosting System: XG Boost

accumulate in a mini-batch manner. When the table has a large number of rows, pre fetching reduces runtime overhead by switching the

direct read/write dependency to a longer dependency.



(a) Allstate 10M



(b) Higgs 10M

Figure9. The approximate algorithm's impact on block size. We discovered that too small blocks result in ineffective parallelization, while too large blocks slow training considerably owing to cache misses.

A non-cache-aware technique was applied on the Higgs and Allstate datasets. The cache-aware implementation of the exact greedy technique runs twice as fast as the naive one when the dataset is large. The problem is handled using approximation methods by choosing the proper block size. The maximum number of examples in a block is determined as the block size. The cost of storing gradient statistics in a cache is reflected in this value. Choosing a block size that is too small reduces

the programmer's workload. As a result, each thread's parallelization is inefficient. On the contrary, Excessively large blocks, on the other hand, cause cache misses. The gradient statistics are too vast for the CPU cache to hold them. This is a good one. The block size balances these two factors. We drew a parallel. Different block sizes were utilised on two data sets. The results Figureshows that balancing the cache property and parallelization by using 216 examples per block

Table 1: Comparison of major tree boosting systems.

System	exact greedy	approximate global	approximate local	out-of-core	sparsity aware	parallel
XGBoost	yes	yes	yes	yes	yes	yes
pGBRT	no	no	yes	no	no	yes
Spark MLlib	no	yes	no	no	partially	yes
H2O	no	yes	no	no	partially	yes
scikit-learn	yes	no	no	no	no	no
R GBM	yes	no	no	no	partially	no

Computation Blocks That Aren't Core

To achieve scalable learning, one of our system's goals is to make the most of a machine's resources. Disk space is necessary in addition to CPUs and memory to handle data that does not fit in main memory. To enable out-of-core processing, we partition the data into many parts and save each block on disc. It's vital to use a separate thread to pre-fetch the block into a main memory buffer during computation so that calculation and disc reading can both happen at the same time. However, because the majority of the computer time is spent reading discs, this does not totally solve the problem. Disc IO overhead must be decreased, and disc IO throughput must be enhanced. We primarily use two ways to improve out-of-core processing.

Blocks are compressed. The first method we use is block compression. When it comes to loading into main memory, there are a few things to keep in We use a general-purpose compression approach to compress the feature values. To retrieve the row index, we subtract the column index from the row index. Use a row index based on the block's beginning index and a block index based on the row index. A 16-bit integer is used to store each offset. There will be 216 examples in total. a single block, which has shown to be a wise decision On the dataset we looked at, we got about a 26 percent to 29 percent accuracy in most cases. the compression ratio It's best to avoid sharding. As a second option, the data can be sharded. onto a variety of dishes in various ways A pre-fetcher is a device that gathers data before sending it to a destination. A thread on each disc collects data and saves it in an array. A memory buffer is a buffer that is retained in memory. The training thread alternately reads the data from each buffer. When a large number of people.

WORKS IN RELATIONSHIP

In our approach, we apply gradient boosting [10] to accomplish additive optimization in functional space. Gradient Tree boosting has been demonstrated to be useful in categorization [12]. Other methods have been tried, such as learning to rank [5], structured prediction [8], and others. XGBoost avoids this by using regularised model overfitting. This is comparable to previous work on regularised greedy forest [25], however the goal and algorithm have been streamlined for

parallelization purposes. Column sampling is a simple but effective sampling method. Random Forest [4] was the source of inspiration for this technique. Sparsity conscious learning is required for other sorts of models, such as financial models. Only a few research on tree learning have looked into linear models [9]. I'd like to discuss this subject in a principled manner. The method proposed in This is the first publication to give a comprehensive approach to dealing with a wide range of issues. sparseness patterns A number of prior works on parallelizing tree learning have been published [22, 19]. The bulk of these algorithms are compatible with the study's paradigm. Data can be partitioned by columns [23], and the precise greedy algorithm can be used. This is also supported by our framework, and techniques like cache-aware pre fetching can aid in the implementation of such an approach. While the vast majority of people Out-of-core computing and cache-aware learning are two undiscovered system directions in which our work advances. Previous research has focused on the algorithmic aspect of parallelization; nevertheless, our findings improve in two previously unknown system directions: out-of-core computation and cache-aware learning. This gives us insight into how the system and algorithm can be enhanced together, resulting in an end-to-end system that can address large-scale problems with a limited amount of computational capacity. We also compared and compare our two businesses. The system and existing open source implementations are shown in Table 1. In the database field, the problem of quantile summary (without weights) is well-known [14, 24]. On the other hand, the approximate tree boosting method demonstrates a larger problem: detecting quantiles on weighted data. To our knowledge, the weighted quantile sketch proposed in this study is the first solution to this problem. The weighted quantile summary isn't just for tree learning; in the future, it could be beneficial in other data science and machine learning applications.

COMPLETE ANALYSIS

The System's Implementation

XGBoost is a free and open source software package⁶. The container is reusable and transportable. It has a user-defined objective function as well as a number of weighted classification and rank objective functions. It's

available in a number of languages, including Python, R, and Julia, and it works in tandem with scikit learn, a language-specific data science library. The all reduction rabbit library⁷ is used in the distributed version. Because XGBoost is mobile, it can be used in a variety of ecosystems rather than being limited to a particular platform. XGBoost is included and works right away. MPI Hadoop and the Sun Grid engine We have added the ability to use XGBoost on jvmbigdata stacks such as Flink

and Spark. The distributed version has also been included. Tianchi⁸ is a cloud-based application. Alibaba is a corporation that Alibaba owns. We are confident that there will be more integrations in the future.

Configuration and Dataset

We employed four different datasets in our research. Table 2 has a brief explanation of these data sets. During a handful of the tests,

Table 2: Dataset used in the Experiments.

Dataset	n	m	Task
Allstate	10M	4227	Insurance claim classification
Higgs Boson	10M	28	Event classification
Yahoo LTRC	473K	700	Learning to Rank
Criteo	1.7B	67	Click through rate prediction

Table 3: Comparison of Exact Greedy Methods with 500 trees on Higgs-1M data.

Method	Time per Tree (sec)	Test AUC
XGBoost	0.6841	0.8304
XGBoost (colsample=0.5)	0.6401	0.8245
scikit-learn	28.51	0.8302
Rgbm	1.032	0.6224

Due to slow baselines or to demonstrate the performance of the algorithm with varied dataset sizes, we use a randomly selected fraction of the data to demonstrate the method's performance with various dataset sizes. We use a suffix to denote the size in various situations. For example, Allstate-10K is a subset of the Allstate dataset that has 10K occurrences.

The first dataset we use is the Allstate insurance claim dataset⁹. The objective is to anticipate the likelihood and cost of an insurance claim based on a variety of risk indicators. By simply forecasting the possibility of an insurance claim, we eased the problem in the experiment. This dataset is used in Section 3.4 to evaluate the sparsity-aware method's impact. The majority of the data's sparse qualities are due to one-hot encoding. The training set is made up of 10M instances at random, while the evaluation set is made up of the rest. The second dataset is the Higgs boson dataset¹⁰ from high-energy physics. The data was generated using Monte Carlo simulations of physical events. It has 21 kinematic properties established by the particle detectors in the accelerator.

It also contains seven other derived physics quantities. of the single particles The purpose

is to determine if an incident is a disaster or not. This symbol is used to represent the Higgs boson. 10M was chosen at random. The remaining examples should be used as a training set, while the rest should be used as an evaluation set. The third dataset is the Yahoo! learning to rank challenge. The dataset [6] is one of the most often used benchmarks in learning to rank algorithms. The information gathered comprises A total of 20K web search requests were made, each of which was associated with a single person. This list contains roughly 22 documents. The purpose is to rate the papers according to how relevant they are to the question. We use the term "official." We had a train-test split in our experiment. The last dataset is the criteo terabyte click log dataset¹¹. This dataset is used to evaluate the system's scalability ability in distributed and out-of-core contexts. The The data contains 13 integer features and 26 user ID features. details about the product and the advertiser Because a tree-based paradigm is built on trees, We preprocess the data to make it more capable of dealing with continuous characteristics. Count statistics and average CTR were used to compile the data. The ID traits were replaced by the appropriate count data for the first ten days, and the ID

A Scalable Tree Boosting System: XG Boost

features were replaced by the equivalent count statistics for the next ten days. a ten-day phase of training After preprocessing, the training set is ready to use. There are 1.7 billion instances with 67 features (13 integer, 13 floating point, 13 floating point, 13.

Average CTR statistics and 26 counts). The entire set of information There is more than one terabyte in Lib SVM format. The first three datasets are for single-machine parallelism, while the fourth is for distributed

and out-of-core computing. A Dell PowerEdge R420 with two eight-core CPUs is used for all single-machine tests. 64GB of RAM with an Intel Xeon (E5-2470) (2.3GHz) processor. If this is the case, all tests are conducted using all of the available resources. There are cores in the machine. The machine settings for distributed and out-of-core tests will be provided in the floating point, 13 floating point, 13.

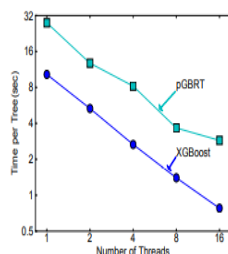


Figure 10: Comparison between XGBoost and pGBRT on Yahoo! LTRC dataset.

Table 4: Comparison of Learning to Rank with 500 trees on Yahoo! LTRC Dataset

Method	Time per Tree (sec)	NDCG@10
XGBoost	0.826	0.7892
XGBoost (colsample=0.5)	0.506	0.7913
pGBRT [22]	2.576	0.7915

Categorization

In this part, we analyse the performance of XGBoost on a single machine using the exact greedy strategy. Higgs-1M data was compared to two additional widely used exact greedy tree boosting techniques. Because scikit-learn only works with non-sparse input, we use it. Use a dense Higgs dataset for a fair comparison. The 1M is the one we use. R's GBM, for example, uses a greedy method to get scikit-learn to finish in a reasonable length of time. XGBoost and scikit-learn both learn the complete tree, however a technique that just develops one branch of a tree makes it

faster but at the sacrifice of precision. The following are the results: The results are shown in Table 3. R's GBM is outperformed by both XGBoost and scikit-learn, with XGBoost requiring longer to run. faster than scikit-learn by a factor of ten We also discovered that column subsamples perform well in this experiment.

Getting a Glimpse of the Ranking Process

After that, XGBoost's performance on the learning to rank task is evaluated. Our findings are compared to those of pGBRT [22]. The best previously released system, XGBoost, was chosen for this competition.

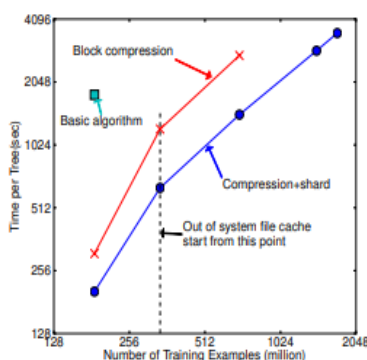


Figure11. Different subsets of criteo data are used to compare out-of-core techniques. Data points are missing due to a shortage of disc space. According to our findings, the simple algorithm can only handle 200 million cases.

A Scalable Tree Boosting System: XG Boost

Compression improves performance by three times, while sharding into two drives improves speed by another two times. The system runs out of file cache after 400M examples. After that, the algorithm will have to rely only on the disc. When the compression + shard technique runs out of file cache, it slows down less severely and then follows a linear trend. While pGBRT supports just an approximation greedy algorithm, runs an accurate greedy algorithm. The results are shown in Table 4 and Figure 10. XG Boost is more effective, as we've discovered. Subsampling columns, it turns out, not only cuts down on running time but also boosts performance for this problem. This could be because subsampling is beneficial.

Try Something New in the Out-of-Core

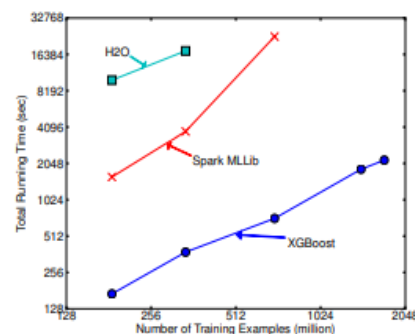
In an out-of-core context, we also use criteo data to examine our system. On a single AWS c3.8xlarge computer, the experiment was carried out (32 cores, two 320 GB SSDs, and 60 GB RAM). The results are shown in Figure 11. Compression accelerates processing by a factor of three, whereas sharding into two discs accelerates computation by a factor of two.

For a true out-of-core scenario in this type of experiment, it's necessary to use a huge dataset

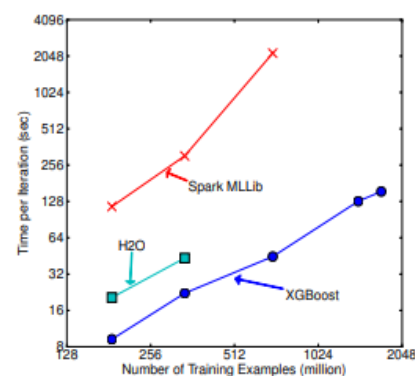
to deplete the system file cache. In reality, this is how we've put things up. We may see a transition point when the system runs out of file cache. It's worth mentioning that the change in the final approach isn't as noticeable. This is due to higher disc throughput and more efficient computational resource utilisation. Our final approach can process 1.7 billion data points on a single system.

Conduct an Experiment with a Diverse Group Of People

Finally, we put the system through its paces in a distributed setting. We used m3 to build up a YARN cluster on EC2. 2xlarge machines, which are frequently used in clusters. Each Eight virtual processors, 30GB of RAM, and two hard drives make up the system. Local SSD discs with an 80GB capacity. The data is stored on AWS S3. to prevent having to pay for long-term storage instead of HDFS We begin by contrasting our system with two systems that are used in production. Two distributed systems are Spark MLlib [18] and H2O 12. We use 32 m3 of space. To test the system's performance, two huge computers are employed.



(a) End-to-end time cost include data loading



(b) Per iteration cost exclude data loading

Figure12. A comparison of multiple distributed systems operating for 10 iterations on a subset of criteo data on 32 EC2 nodes.

A Scalable Tree Boosting System: XG Boost

XGBoost is 10 times quicker than Spark and 2.2 times faster than H2O's optimised version every iteration (although, H2O is slow in loading data, resulting in a longer end-to-end time). It's worth noting that spark has been affected. You'll notice a substantial slowdown when you run out of memory. XGBoost is a more efficient and scalable alternative. With the resources offered by Out-of-core computation, the complete 1.7 billion examples are exploited. variables with varying input sizes Both of the baseline systems are in-memory analytics frameworks that require data to be kept in RAM; however, once memory is exhausted, XGBoost can switch to an out-of-core mode. The results are on show. We discovered that XGBoost surpasses the competition in Figure 12.

We discovered that XGBoost surpasses the competition in Figure 12. The underlying

systems More importantly, it is capable of taking on new tasks.

Scale up seamlessly by utilising out-of-core computing.

We were able to complete all 1.7 billion cases despite the restricted processing resources available. Only a portion of the data can be handled by the baseline systems of the data with the resources at hand This experiment proves that the benefit of integrating all system changes together and solving a large-scale real-world problem We also evaluate The scaling feature of XGBoost can be modified by adjusting the number of iterations machines. The results are shown in Figure 13. We have the ability to locate XGBoost's performance scales linearly as we add more machines.

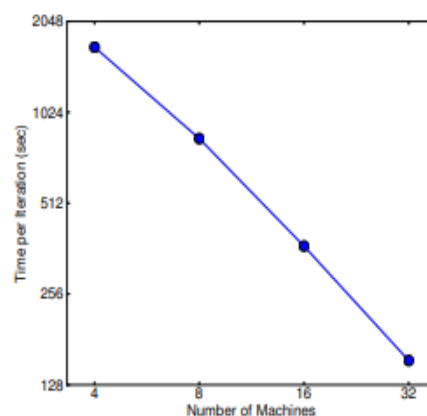


Figure13. XGBoost scaling on the Criteo entire 1.7 billion dataset with varied numbers of computers. Using more workstations expands the file cache and speeds up the system, perpetuating the trend. to have a modest non-linearity XGBoost can handle the entire dataset on as few as four workstations, and it expands gracefully when more resources become available.

FINAL COMMENTS

This paper discusses the lessons we learnt while designing XGBoost, a scalable tree boosting method. It appeals to data scientists because it provides them with access to cutting-edge technologies. On a range of topics, the truth comes out. A new sparsity model was proposed. a theoretically sound sparse data processing approach A justified weighted quantile sketch is used for approximate learning. According to our experience, cache access patterns, data compression, and sharding are essential components for building a high-performance database. A scalable end-to-end system is required for tree boosting.

ACKNOWLEDGMENTS

We received helpful feedback from Tyler B. Johnson, Marco Tulio Ribeiro, Sameer Singh, and Arvind Krishnamurthy. We'd also like to thank Tong He, Bing Xu, Michael Benesty, Yuan Tang, Hongliang Liu, Qiang Kou, Nan Zhu, and the rest of the XGBoost community. This effort was funded by the Office of Naval Research (N000141010672), the National Science Foundation (NSF IIS 1258741), and the MARCO-sponsored Terra Swarm Research Center DARPA.

REFERENCES

- [1] G. Ridgeway. Generalized Boosted Models: A guide to the gbm package.
- [2] L. Breiman. Random forests. *Maching Learning*, 45(1):5–32, Oct. 2001.
- [3] M. Greenwald and S. Khanna. Space-efficient online computation of quantile summaries. In *Proceedings of the 2001 ACM SIGMOD International Conference on Management of Data*, pages 58–66, 2001.
- [4] J. Ye, J.-H. Chow, J. Chen, and Z. Zheng. Stochastic gradient boosted distributed decision trees. In *Proceedings of the 18th ACM Conference on Information and Knowledge Management, CIKM '09*.
- [5] P. Li. Robust Logitboost and adaptive base class (ABC) Logitboost. In *Proceedings of the Twenty-Sixth Conference Annual Conference on Uncertainty in Artificial Intelligence (UAI'10)*, pages 302–311, 2010.
- [6] S. Tyree, K. Weinberger, K. Agrawal, and J. Paykin. Parallel boosted regression trees for web search ranking. In *Proceedings of the 20th international conference on World wide web*, pages 387–396. ACM, 2011.
- [7] J. Friedman. Greedy function approximation: a gradient boosting machine. *Annals of Statistics*, 29(5):1189–1232, 2001.
- [8] J. H. Friedman and B. E. Popescu. Importance sampled learning ensembles, 2003.
- [9] P. Li, Q. Wu, and C. J. Burges. Mcrank: Learning to rank using multiple classification and gradient boosting. In *Advances in Neural Information Processing Systems 20*, pages 897–904. 2008.
- [10] R.-E. Fan, K.-W. Chang, C.-J. Hsieh, X.-R. Wang, and C.-J. Lin. LIBLINEAR: A library for large linear classification. *Journal of Machine Learning Research*, 9:1871–1874, 2008.
- [11] J. Friedman. Stochastic gradient boosting. *Computational Statistics & Data Analysis*, 38(4):367–378, 2002.
- [12] X. He, J. Pan, O. Jin, T. Xu, B. Liu, T. Xu, Y. Shi, A. Atallah, R. Herbrich, S. Bowers, and J. Q. n. Candela. Practical lessons from predicting clicks on ads at facebook. In *Proceedings of the Eighth International Workshop on Data Mining for Online Advertising, ADKDD'14*, 2014.
- [13] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.
- [14] T. Chen, S. Singh, B. Taskar, and C. Guestrin. Efficient second-order gradient boosting for conditional random fields. In *Proceeding of 18th Artificial Intelligence and Statistics Conference (AISTATS'15)*, volume 1, 2015.
- [15] T. Zhang and R. Johnson. Learning nonlinear functions using regularized greedy forest. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 36(5), 2014.
- [16] B. Panda, J. S. Herbach, S. Basu, and R. J. Bayardo. Planet: Massively parallel learning of tree ensembles with mapreduce. *Proceeding of VLDB Endowment*, 2(2):1426–1437, Aug. 2009.
- [17] X. Meng, J. Bradley, B. Yavuz, E. Sparks, S. Venkataraman, D. Liu, J. Freeman, D. Tsai, M. Amde, S. Owen, D. Xin, R. Xin, M. J. Franklin, R. Zadeh, M. Zaharia, and A. Talwalkar. MLlib: Machine learning in apache spark. *Journal of Machine Learning Research*, 17(34):1–7, 2016.
- [18] Q. Zhang and W. Wang. A fast algorithm for approximate quantiles in high speed data streams. In *Proceedings of the 19th International Conference on Scientific and Statistical Database Management*, 2007.
- [19] X. Meng, J. Bradley, B. Yavuz, E. Sparks, S. Venkataraman, D. Liu, J. Freeman, D. Tsai, M. Amde, S. Owen, D. Xin, R. Xin, M. J. Franklin, R. Zadeh, M. Zaharia, and A. Talwalkar. MLlib: Machine learning in apache spark. *Journal of Machine Learning Research*, 17(34):1–7, 2016.
- [20] Q. Zhang and W. Wang. A fast algorithm for approximate quantiles in high speed data streams. In *Proceedings of the 19th International Conference on Scientific and Statistical Database Management*, 2007.
- [21] O. Chapelle and Y. Chang. Yahoo! Learning to Rank Challenge Overview. *Journal of Machine Learning Research - W & CP*, 14:1–24, 2011.
- [22] J. Friedman, T. Hastie, and R. Tibshirani. Additive logistic regression: a statistical view of boosting. *Annals of Statistics*, 28(2):337–407, 2000.
- [23] X. He, J. Pan, O. Jin, T. Xu, B. Liu, T. Xu, Y. Shi, A. Atallah, R. Herbrich, S. Bowers, and J. Q. n. Candela. Practical lessons from predicting clicks on ads at facebook. In *Proceedings of the Eighth International Workshop on Data Mining for Online Advertising, ADKDD'14*, 2014.

- [24] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.
- [25] T. Chen, S. Singh, B. Taskar, and C. Guestrin. Efficient second-order gradient boosting for conditional random fields. In *Proceeding of 18th Artificial Intelligence and Statistics Conference (AISTATS'15)*, volume 1, 2015.

Citation: Mounika Nalluri et al, “A Scalable Tree Boosting System: XG Boost”, *International Journal of Research Studies in Science, Engineering and Technology* 2020; 7(12): 36-51. DOI: <https://doi.org/10.22259/2349-476X.0712005>

Copyright: © 2020 Authors. This is an open-access article distributed under the terms of the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

