# XGBoost: A Scalable Tree Boosting System
## A Review and Analysis

Adrien Protzel

AI 534 - Machine Learning

Instructor: Liang Huang

The paper "XGBoost: A Scalable Tree Boosting System" addresses the central problem of scaling gradient tree boosting to large, sparse datasets that often exceed main memory. Gradient boosting is a leading method for classification and ranking in production settings such as click-through rate prediction, search, and fraud detection, yet traditional implementations incur heavy sorting and memory bandwidth costs, are sensitive to cache behavior, and do not natively handle missing values or weighted approximate split proposals. XGBoost addresses these issues with a careful combination of algorithmic advances and systems engineering that together make tree boosting efficient on single machines, with out-of-core data streams, and across distributed clusters.

The work was authored by Tianqi Chen and Carlos Guestrin at the University of Washington, with Guestrin serving as principal investigator and Chen as the student author. It was published at KDD 2016 in San Francisco, and the canonical version is the conference proceedings paper; there is no widely cited journal extension beyond the arXiv preprint. The open source implementation became a standard tool for tabular machine learning with bindings across Python, R, JVM, and C++, and it has a large contributor community and broad adoption in competitions and industry.
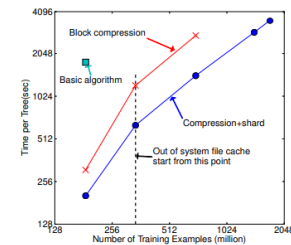


Figure 11: Comparison of out-of-core methods on different subsets of criteo data. The missing data points are due to out of disk space. We can find that basic algorithm can only handle 200M examples. Adding compression gives 3x speedup, and sharding into two disks gives another 2x speedup. The system runs out of file cache start from 400M examples. The algorithm really has to rely on disk after this point. The compression+shard method has a less dramatic slowdown when running out of file cache, and exhibits a linear trend afterwards.

As a non-technical background, the paper was presented at KDD in August 2016. There is no record of a formal best-paper award specific to this submission, but the system's community impact is clear. The authors note that in 2015, among 29 challenge-winning solutions on Kaggle, 17 used XGBoost, and it was used by every winning team in the top ten of KDD Cup 2015. For impact, bibliometric aggregators show very high citation counts; for example, SciSpace reports roughly fifteen thousand citations as of late 2025, reflecting sustained academic and practical influence over time. Talks and media are available, including a KDD 2016 presentation video and community write-ups and workshops that amplified the system's reach.

On reproducibility, the paper did not introduce a new dataset. It used Allstate insurance claims, the Higgs Boson dataset, the Yahoo Learning to Rank Challenge data, and the Criteo terabyte click logs. The Yahoo LTRC dataset is accessible under specific terms and includes an official train, validation, and test split; public catalog pages document the split sizes and gated access requirements. The Criteo terabyte logs are publicly hosted and widely used for CTR benchmarks, with documentation on structure and download sources. The authors released code as an open source package with language bindings and distributed training support, and community demo notebooks and examples exist in the XGBoost repository and broader ecosystem, making hands-on reproduction practical. Slide decks and talk videos are available from KDD and workshops, offering additional implementation detail. Overall, the paper is reproducible for most experiments if datasets can be obtained under their terms and if hardware and configuration are carefully matched.

The learning objective combines a differentiable loss with an explicit regularizer over the trees to control complexity: [Section 2.1, Eq. (2)]

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

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

The model is an additive ensemble of trees that sums the contribution of each tree to make a prediction: [Section 2.1, Eq. (1)]

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

Training proceeds in additive steps using a second order approximation of the loss around current predictions, with per instance gradient and Hessian statistics guiding the update: [Section 2.2, Eq. (3)]

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

For a fixed tree structure, the algorithm uses a closed form expression for the optimal leaf weight: [Section 2.2, Eq. (5)]

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

It evaluates candidate splits with a gain score that compares the children to the parent: [Section 2.2, Eq. (7)]

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

These analytic scores make tree growth and pruning efficient and well regularized on large, sparse tabular data.

---

**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**
    $G_L \leftarrow G_L + g_j, H_L \leftarrow H_L + h_j$
    $G_R \leftarrow G - G_L, H_R \leftarrow H - H_L$
    $score \leftarrow max(score, \frac{G_L^2}{H_L+\lambda} + \frac{G_R^2}{H_R+\lambda} - \frac{G^2}{H+\lambda})$
  **end**
**end**
**Output:** Split with max score

---

**Algorithm 2:** Approximate Algorithm for Split Finding

**for** $k = 1$ **to** $m$ **do**
  Propose $S_k = \{s_{k1}, s_{k2}, \cdots s_{kl}\}$ by percentiles on feature $k$.
  Proposal can be done per tree (global), or per split(local).
**end**
**for** $k = 1$ **to** $m$ **do**
  $G_{kv} \leftarrow = \sum_{j\in\{j|s_{k,v}\geq \mathbf{x}_{jk} > s_{k,v-1}\}} g_j$
  $H_{kv} \leftarrow = \sum_{j\in\{j|s_{k,v}\geq \mathbf{x}_{jk} > s_{k,v-1}\}} h_j$
**end**
Follow same step as in previous section to find max score only among proposed splits.

---

What the paper solves is scaling gradient tree boosting across single-machine, out-of-core, and distributed settings while addressing sparsity, cache behavior, and weighted approximate split proposals. Prior systems were fast but constrained by sorting and memory bandwidth, handled missing values inconsistently, and lacked a correct weighted quantile method. Why this matters is that boosted trees are a dominant choice for tabular prediction in production, so solving scalability and sparsity unlocks practical accuracy and speed on massive logs and sparse features. How it works is a mix of algorithmic and system design: a weighted quantile sketch for approximate

split proposals; sparsity-aware split finding with default directions for missing values; a compressed, sorted column-block layout; cache-aware prefetching; out-of-core training with compression and sharding; and a distributed allreduce interface. The wow factor is strong empirical performance: on Higgs-1M, XGBoost achieved competitive AUC with dramatically lower time per tree than scikit-learn, and on Yahoo LTRC it matched or exceeded ranking metrics while training faster; the system also scaled to 1.7 billion examples on a single machine and across clusters.

The system design is as important as the algorithms. XGBoost organizes data by feature in compressed, sorted column blocks so the learner can scan each column once and evaluate splits without repeatedly sorting. It keeps computations close to the CPU through cache-aware prefetching, which reduces stalls and helps the training loop move smoothly. When the dataset is larger than memory, XGBoost streams blocks from disk, compresses them to reduce I/O, and shards them across multiple disks so reading and computation can proceed in parallel. In a cluster, it coordinates workers with reliable collective communication, which lets the same training procedure scale from a single machine to many machines with only minor changes.
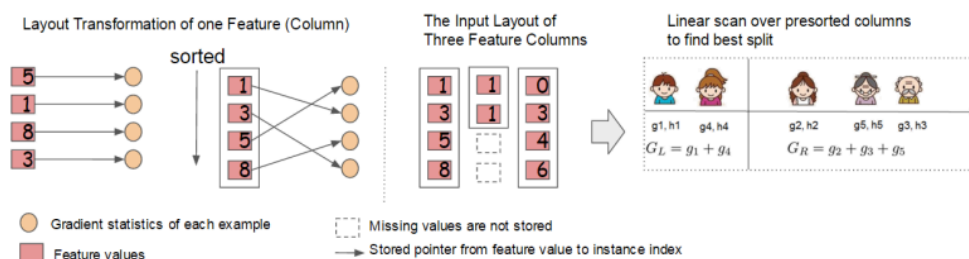


Figure 6: Block structure for parallel learning. Each column in a block is sorted by the corresponding feature value. A linear scan over one column in the block is sufficient to enumerate all the split points.

In practice, these choices make XGBoost feel fast and resilient. On common classification and ranking tasks it reaches strong accuracy while building trees quickly. The sparsity-aware split finder focuses work on non-missing values and learns a sensible default path for missing entries, which is helpful for one-hot and sparse features. The cache-friendly layout and the streaming approach allow large log data to be trained without needing all of it in memory, and the distributed version extends the same workflow across a small cluster when more capacity is needed.

There are additional drawbacks to consider beyond documentation and replicability. Because the supplementary proofs for the weighted quantile sketch have not always been easy to access, some readers find it harder to validate the theory behind the approximate split proposals, and obtaining datasets like Yahoo LTRC under their terms adds friction to reproduction. Even when the data are available, differences in library versions, default settings, and hardware can lead to small variations in results, so reproducing benchmarks exactly may require careful control of seeds, environments, and configuration. Out-of-core training depends on disk throughput and storage layout, so performance can vary across systems, and cluster training introduces practical concerns around networking, fault tolerance, and resource management that are outside the scope of the core algorithm. Model transparency is limited compared to single decision trees, and many deployments add post-hoc explanation, calibration, and fairness checks, which increase operational work. Hyperparameter tuning can be

substantial, and handling very high-cardinality categoricals often requires careful preprocessing.

Regarding relevance to the course, the perceptron and other linear models gave me a baseline for thinking about simple decision boundaries, and reading this paper showed how boosted trees extend that idea by layering many small corrections to handle nonlinear structure. Our work with KNN highlighted instance-based learning and the importance of feature representation; XGBoost's ability to focus on non-missing values and treat one-hot zeros sensibly connects directly to how KNN behaved on sparse inputs. In the price prediction exercises, we saw how regularization and careful feature use can stabilize regression, and XGBoost's shrinkage and tree complexity controls mirror those lessons while letting the model capture interactions we could not express linearly. In the review categorization problem, we dealt with high-cardinality, mostly sparse text features and occasional missing signals, and XGBoost's sparsity-aware split finding and learned default paths fit that setting well. Finally, the paper's attention to practical system details such as column-wise storage, cache-friendly access, and streaming data from disk helped explain the training bottlenecks we noticed whenever data or features grew beyond memory in our labs.