

Extreme Gradient Boosting (XGBoost)

Luu Minh Sao Khue

December 8, 2025

1. Key Ideas and Comparison to Gradient Boosting Trees

Gradient Boosting Trees (GBT) build a strong predictor by adding many weak regression trees in a stage-wise manner. At each boosting iteration, a new tree is fitted to the residual errors (or gradients) of the current model, and the final prediction is the sum of all tree outputs.

XGBoost (Extreme Gradient Boosting) follows the same general idea but improves it in three major ways:

- **Smarter math:** It uses both first- and second-order derivatives of the loss function with respect to the predictions, allowing Newton-style updates of tree leaf values and principled split scoring.
- **Stronger regularization:** It includes explicit regularization on tree complexity (number of leaves) and on leaf weights, reducing overfitting and encouraging simpler trees.
- **Efficient implementation:** It incorporates engineering tricks such as histogram-based split finding, column and row subsampling, and sparse-aware computation, which make it very fast and scalable.

Compared to a standard GBT implementation that typically uses only first-order information (gradients) and heuristic leaf values, XGBoost optimizes a regularized objective that directly includes second-order information and complexity penalties. This leads to more accurate trees, better generalization, and faster convergence.

2. Intuition Behind XGBoost

Consider a dataset with n samples. For each sample i , we have:

- feature vector $x_i \in \mathbb{R}^m$,
- target y_i ,
- model prediction \hat{y}_i .

Instead of learning a single large model, gradient boosting builds the prediction function as a sum of many small regression trees:

$$\hat{y}_i^{(T)} = \sum_{t=1}^T f_t(x_i), \quad (1)$$

where f_t is the regression tree added at boosting iteration t . Each tree f_t is trained to correct the mistakes of the current model $\hat{y}^{(t-1)}$.

XGBoost improves this process by:

- using the first derivative (gradient) to know in which direction to correct the prediction,
- using the second derivative (Hessian) to know how confidently to move in that direction (how big the correction should be),
- choosing tree splits by how much they reduce a regularized approximation of the loss,
- computing the best leaf values in closed form, instead of by heuristic.

Thus each new tree is chosen to give the largest reduction in a regularized second-order Taylor approximation of the overall loss. The result is a sequence of small, carefully chosen corrections that gradually shape a strong and well-regularized predictor.

3. Mathematical Formulation

In this section we derive the XGBoost formulation step by step and explain all notation.

3.0.1 Notation and Regularized Objective

We have:

- n : number of training samples,
- m : number of features,
- T : number of boosting iterations (trees),
- $x_i \in \mathbb{R}^m$: feature vector of sample i ,
- y_i : target/label of sample i ,
- $\hat{y}_i^{(t)}$: prediction for sample i after t trees,
- f_t : regression tree added at iteration t .

The model after T trees is

$$\hat{y}_i^{(T)} = \sum_{t=1}^T f_t(x_i). \quad (2)$$

We define a regularized objective:

$$\mathcal{L} = \sum_{i=1}^n l(y_i, \hat{y}_i^{(T)}) + \sum_{t=1}^T \Omega(f_t), \quad (3)$$

where:

- $l(y_i, \hat{y}_i)$ is the loss for one sample (e.g., squared error or logistic loss),
- $\Omega(f_t)$ is the complexity penalty for tree f_t .

In XGBoost, a tree f is regularized as:

$$\Omega(f) = \gamma K + \frac{1}{2} \lambda \sum_{j=1}^K w_j^2, \quad (4)$$

where:

- K is the number of leaves (terminal nodes) in the tree,
- $w_j \in \mathbb{R}$ is the prediction value associated with leaf j ,
- $\gamma \geq 0$ is a penalty per leaf (discourages too many leaves),
- $\lambda \geq 0$ is the L2 regularization on leaf values.

Improvement: Including γ and λ in the objective ensures that tree complexity is controlled during training, reducing overfitting and encouraging smaller, simpler trees.

3.1 Boosting Iteration and Second-Order Taylor Approximation

Boosting proceeds by adding trees one by one. Suppose we have already built trees f_1, \dots, f_{t-1} , giving predictions

$$\hat{y}_i^{(t-1)} = \sum_{k=1}^{t-1} f_k(x_i). \quad (5)$$

We now add a new tree f_t :

$$\hat{y}_i^{(t)} = \hat{y}_i^{(t-1)} + f_t(x_i). \quad (6)$$

The objective at iteration t can be written as

$$\mathcal{L}^{(t)} = \sum_{i=1}^n l(y_i, \hat{y}_i^{(t-1)} + f_t(x_i)) + \Omega(f_t) + (\text{constant terms from previous trees}). \quad (7)$$

The terms from previous trees do not depend on f_t and can be ignored when optimizing f_t . We focus on the effective objective:

$$\tilde{\mathcal{L}}^{(t)} = \sum_{i=1}^n l(y_i, \hat{y}_i^{(t-1)} + f_t(x_i)) + \Omega(f_t). \quad (8)$$

We approximate the loss using a second-order Taylor expansion with respect to the prediction. Let

$$g_i = \frac{\partial l(y_i, \hat{y})}{\partial \hat{y}} \Big|_{\hat{y}=\hat{y}_i^{(t-1)}}, \quad h_i = \frac{\partial^2 l(y_i, \hat{y})}{\partial \hat{y}^2} \Big|_{\hat{y}=\hat{y}_i^{(t-1)}}, \quad (9)$$

where:

- g_i is the first derivative of the loss with respect to the prediction (gradient),
- h_i is the second derivative of the loss with respect to the prediction (Hessian).

For each sample i we approximate:

$$l(y_i, \hat{y}_i^{(t-1)} + f_t(x_i)) \approx l(y_i, \hat{y}_i^{(t-1)}) + g_i f_t(x_i) + \frac{1}{2} h_i f_t(x_i)^2. \quad (10)$$

Summing over i and dropping the constant term $\sum_i l(y_i, \hat{y}_i^{(t-1)})$ we obtain the approximate objective:

$$\hat{\mathcal{L}}^{(t)} = \sum_{i=1}^n \left[g_i f_t(x_i) + \frac{1}{2} h_i f_t(x_i)^2 \right] + \Omega(f_t). \quad (11)$$

Improvement: Using a second-order Taylor approximation allows XGBoost to include curvature information (via h_i) and to turn the problem of choosing leaf values into a simple quadratic optimization.

3.2 Tree Structure and Leaf Statistics

Each tree f_t is a piecewise constant function:

$$f_t(x) = w_{q(x)}, \quad (12)$$

where:

- $q(x)$ is the structure of the tree, mapping an input x to a leaf index $j \in \{1, \dots, K\}$,
- w_j is the prediction value of leaf j ,
- K is the number of leaves.

Define the index set of samples that fall into leaf j :

$$I_j = \{i \mid q(x_i) = j\}. \quad (13)$$

Then for $i \in I_j$ we have $f_t(x_i) = w_j$, and

$$\sum_{i=1}^n g_i f_t(x_i) = \sum_{i=1}^n g_i w_{q(x_i)} = \sum_{j=1}^K \left(\sum_{i \in I_j} g_i \right) w_j, \quad (14)$$

$$\sum_{i=1}^n \frac{1}{2} h_i f_t(x_i)^2 = \sum_{i=1}^n \frac{1}{2} h_i w_{q(x_i)}^2 = \sum_{j=1}^K \left(\sum_{i \in I_j} \frac{1}{2} h_i \right) w_j^2. \quad (15)$$

Define the aggregated gradient and Hessian for leaf j as:

$$G_j = \sum_{i \in I_j} g_i, \quad H_j = \sum_{i \in I_j} h_i. \quad (16)$$

Plugging these into the objective:

$$\hat{\mathcal{L}}^{(t)} = \sum_{j=1}^K \left[G_j w_j + \frac{1}{2} H_j w_j^2 \right] + \gamma K + \frac{1}{2} \lambda \sum_{j=1}^K w_j^2. \quad (17)$$

We can group terms for each leaf:

$$\hat{\mathcal{L}}^{(t)} = \sum_{j=1}^K \left[G_j w_j + \frac{1}{2} (H_j + \lambda) w_j^2 \right] + \gamma K. \quad (18)$$

Improvement: Summarizing each leaf by just (G_j, H_j) allows efficient computation of the objective and its minimizer, and enables fast split evaluation during tree construction.

3.3 Optimal Leaf Weights

For a fixed tree structure (i.e., fixed sets I_j), we optimize over leaf weights w_j . For each leaf j we consider:

$$\phi_j(w_j) = G_j w_j + \frac{1}{2} (H_j + \lambda) w_j^2. \quad (19)$$

Differentiating with respect to w_j and setting to zero:

$$\frac{d\phi_j}{dw_j} = G_j + (H_j + \lambda) w_j = 0, \quad (20)$$

gives the optimal leaf weight

$$w_j^* = -\frac{G_j}{H_j + \lambda}. \quad (21)$$

Interpretation:

- G_j is the total gradient in the leaf; $-G_j$ points in the direction that reduces the loss.
- $H_j + \lambda$ in the denominator scales the update by the curvature plus regularization. Larger curvature H_j or larger λ make the update smaller, stabilizing the model and shrinking leaf values.

Improvement: This closed-form solution provides an exact optimal update for leaf values under the quadratic approximation, in contrast to heuristic updates in plain GBT.

3.4 Score of a Tree Structure

Substituting w_j^* into $\hat{\mathcal{L}}^{(t)}$, we obtain the minimal value of the approximate objective for a given tree structure q :

$$\hat{\mathcal{L}}^{(t)}(q) = -\frac{1}{2} \sum_{j=1}^K \frac{G_j^2}{H_j + \lambda} + \gamma K. \quad (22)$$

This quantity (up to a constant) is used as the *score* of the tree structure. Lower values correspond to better (lower loss) trees.

3.5 Split Gain Formula

When growing a tree, XGBoost decides whether or not to split a node by computing the *gain* in the objective. Consider a node containing index set I , with:

$$G = \sum_{i \in I} g_i, \quad H = \sum_{i \in I} h_i. \quad (23)$$

If this node is a leaf, its contribution to the objective is

$$\text{Score}_{\text{parent}} = -\frac{1}{2} \frac{G^2}{H + \lambda} + \gamma. \quad (24)$$

Suppose we split this node into left and right children with index sets I_L and I_R . Define:

$$G_L = \sum_{i \in I_L} g_i, \quad H_L = \sum_{i \in I_L} h_i, \quad G_R = \sum_{i \in I_R} g_i, \quad H_R = \sum_{i \in I_R} h_i. \quad (25)$$

Then the contribution of the two children is

$$\text{Score}_{\text{children}} = -\frac{1}{2} \frac{G_L^2}{H_L + \lambda} - \frac{1}{2} \frac{G_R^2}{H_R + \lambda} + 2\gamma. \quad (26)$$

The gain from performing this split is defined as

$$\text{Gain} = \text{Score}_{\text{children}} - \text{Score}_{\text{parent}}. \quad (27)$$

After simplification this gives the standard XGBoost gain formula:

$$\text{Gain} = \frac{1}{2} \left(\frac{G_L^2}{H_L + \lambda} + \frac{G_R^2}{H_R + \lambda} - \frac{G^2}{H + \lambda} \right) - \gamma. \quad (28)$$

A split is accepted if Gain exceeds zero (or exceeds a user-specified minimum gain).

Improvement: The gain formula combines gradient and curvature information with regularization:

- The terms $\frac{G_L^2}{H_L + \lambda}$ and $\frac{G_R^2}{H_R + \lambda}$ reflect how much the loss can be reduced in each child, scaled by curvature and regularization.
- The term $\frac{G^2}{H + \lambda}$ represents the score of the unsplit node.
- The penalty γ ensures that small improvements do not lead to extra splits.

Thus XGBoost prefers splits that significantly reduce the regularized approximate loss and automatically controls tree complexity.

3.6 Examples of Gradients and Hessians

For completeness, we give two common examples of g_i and h_i .

Squared error regression. For loss

$$l(y_i, \hat{y}_i) = \frac{1}{2}(y_i - \hat{y}_i)^2, \quad (29)$$

we have

$$g_i = \frac{\partial l}{\partial \hat{y}_i} = \hat{y}_i - y_i, \quad h_i = \frac{\partial^2 l}{\partial \hat{y}_i^2} = 1. \quad (30)$$

Binary logistic loss. Let \hat{y}_i be the logit and

$$p_i = \sigma(\hat{y}_i) = \frac{1}{1 + e^{-\hat{y}_i}}, \quad y_i \in \{0, 1\}. \quad (31)$$

The loss is

$$l(y_i, \hat{y}_i) = -[y_i \log p_i + (1 - y_i) \log(1 - p_i)]. \quad (32)$$

Then

$$g_i = \frac{\partial l}{\partial \hat{y}_i} = p_i - y_i, \quad h_i = \frac{\partial^2 l}{\partial \hat{y}_i^2} = p_i(1 - p_i). \quad (33)$$

These derivatives are plugged into the general framework above to compute G_j , H_j , optimal w_j^* , and the split gains.

4. Engineering Optimizations in XGBoost

In addition to its mathematical improvements, XGBoost derives much of its practical effectiveness from several system-level optimizations that greatly accelerate training. These engineering techniques make XGBoost scalable to large datasets and robust in real-world settings.

Histogram-based split finding. Instead of scanning every unique feature value when searching for the best split, XGBoost groups continuous features into histogram bins. Statistics are accumulated per bin, and splits are evaluated using these summaries. This reduces computation, improves memory locality, and allows multi-threaded parallel processing of features.

Column and row subsampling. XGBoost allows random subsampling of both features and data rows for each tree or even each split. This reduces overfitting, accelerates training, and increases model diversity—similar to the benefits seen in Random Forests.

Sparse-aware learning. Real-world datasets often contain missing or zero values. XGBoost treats missing values as informative and learns an optimal default direction (left or right) for them during split construction. This avoids manual imputation and makes the algorithm naturally efficient on sparse matrices.

Cache-aware and block-structured computation. Data are stored in a compressed column format with blocks arranged to maximize cache efficiency. This improves memory access patterns during split evaluation and yields substantial speedups, especially on large datasets with many features.

Out-of-core training. When the dataset does not fit into memory, XGBoost supports out-of-core training by using on-disk storage with prefetching and data compression. This enables the model to scale to very large datasets without requiring specialized hardware.

These engineering optimizations complement the second-order optimization and regularization framework, making XGBoost both highly accurate and extremely efficient in practice.

5. Comparison Between Gradient Boosting Trees and XGBoost

GBT and XGBoost follow the same boosting principle—iteratively adding decision trees to correct the residual errors of previous trees. However, XGBoost introduces several theoretical and engineering improvements that substantially enhance accuracy, regularization, and computational efficiency. Table 1 summarizes the key differences.

5.1 Optimization Formulation

GBT: Classical Gradient Boosting Trees use a first-order approximation of the loss function. At iteration m , the model is updated by fitting a tree to the negative gradients (also called *pseudo-residuals*):

$$r_i^{(m)} = -\frac{\partial \ell(y_i, \hat{y}_i^{(m-1)})}{\partial \hat{y}_i}.$$

This quantity represents the direction in which the model’s prediction should change to reduce the loss most rapidly. The approximation is purely linear, as GBT does not use second-order (curvature) information.

XGBoost: XGBoost generalizes GBT by using a second-order Taylor expansion of the loss around the current prediction:

$$\ell(y_i, \hat{y}_i + f(x_i)) \approx \ell(y_i, \hat{y}_i) + g_i f(x_i) + \frac{1}{2} h_i f(x_i)^2,$$

where $g_i = \frac{\partial \ell}{\partial \hat{y}_i}$ and $h_i = \frac{\partial^2 \ell}{\partial \hat{y}_i^2}$. The term $g_i f(x_i)$ plays the same role as the first-order correction used in GBT (since GBT fits a tree to approximate $-g_i$), while the additional quadratic term $\frac{1}{2} h_i f(x_i)^2$ incorporates curvature and stabilizes the update. Minimizing this second-order approximation yields a closed-form optimal leaf weight:

$$w^* = -\frac{\sum_i g_i}{\sum_i h_i + \lambda},$$

which integrates both gradient and curvature information as well as L2 regularization.

5.2 Leaf Value Computation

GBT: Leaf predictions are typically the mean negative residual in each leaf (or the minimizer of the loss restricted to that leaf). This uses only first-order information and does not incorporate explicit regularization.

XGBoost: Each leaf weight is determined by the closed-form optimum of the regularized second-order objective:

$$w^* = -\frac{G}{H + \lambda},$$

where $G = \sum_i g_i$ and $H = \sum_i h_i$. As a result, leaf values adapt to both the gradient and curvature of the loss, with L2 regularization applied automatically.

5.3 Split Selection

GBT: Splits are selected by maximizing the reduction in residual error or improvement in the first-order approximation of the loss. There is no explicit penalty for model complexity, and the criterion does not use second-order information.

XGBoost: XGBoost uses a regularized *gain* function:

$$\text{Gain} = \frac{1}{2} \left(\frac{G_L^2}{H_L + \lambda} + \frac{G_R^2}{H_R + \lambda} - \frac{G^2}{H + \lambda} \right) - \gamma,$$

which accounts for both the first- and second-order statistics and penalizes unnecessary splits using the regularization terms λ and γ . This leads to more stable and conservative split decisions.

5.4 Regularization

GBT: Regularization in classical GBT is primarily heuristic or structural. Common techniques include shrinkage (learning rate), limiting tree depth, and subsampling of rows or columns. These methods help reduce overfitting but do not appear explicitly in the objective function.

XGBoost: XGBoost introduces explicit regularization in the objective:

$$\Omega(f) = \gamma T + \frac{1}{2} \lambda \sum_{j=1}^T w_j^2,$$

where T is the number of leaves. This penalizes both large leaf weights and overly complex trees, providing principled control over model capacity.

5.5 Engineering and Computational Efficiency

GBT: Traditional GBT implementations rely on exact greedy splitting and do not exploit parallelism effectively. They may struggle with sparse features and large datasets.

XGBoost: XGBoost introduces several engineering optimizations: histogram-based split search, parallel computation, sparsity-aware algorithms, cache optimization, and out-of-core training for large datasets. These enhancements significantly improve scalability and speed.

Table 1: Summary of Key Differences Between GBT and XGBoost

Aspect	GBT	XGBoost
Optimization	First-order (gradient only)	Second-order (gradient + Hessian)
Leaf value computation	Mean residual / first-order minimizer	Closed-form $-G/(H + \lambda)$
Split selection	Residual error reduction	Regularized gain with λ, γ
Regularization	Shrinkage, depth limits, subsampling	Explicit objective regularization
Engineering	Basic tree building	Parallel, histogram-based, sparse-aware