# XGBoost: a Scalable Tree Boosting System.

Tianqi Chen, Carlos Guestrin (2016)

## Gradient tree boosting

Let a dataset with $n$ examples and $m$ features be denoted by

$$D = \{(x_i, y_i)\} \; (|D| = n, \; x_i \in \mathbb{R}^m, \; y_i \in \mathbb{R})$$

A tree ensemble with $K$ additive functions is used to predict output

$$\hat{y}_i = \phi(x_i) = \sum_{k=1}^{K} f_k(x_i) \;, \quad f_k \in \mathcal{F} \rightsquigarrow \text{space of trees}$$

$$\mathcal{F} = \{f(x) = w_{q(x)}\} \; (q: \mathbb{R}^m \to T, \; w \in \mathbb{R}^T)$$

maps example to leaf index    number of leaves    leaf weights

To learn the functions, we use the regularised objective:

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

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

Because this is a difficult optimisation, we use the greedy additive method.

Taylor expansion $\Bigg($

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

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

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

$h_i = \partial^2_{\hat{y}^{(t-1)}} l(y_i, \hat{y}^{(t-1)})$

· Let $I_j = \{ i \mid q(x_i) = j \}$ be the instance set of leaf $j$ (ie the set containing all examples in that leaf.

Removing constant terms and regrouping gives:

$$\mathcal{L}^{(t)} = \sum_{j=1}^{T} \left[ \left( \underbrace{\sum_{i \in I_j} g_i}_{G_j} \right) w_j + \frac{1}{2} \left( \underbrace{\sum_{i \in I_j} h_i + \lambda}_{H_j} \right) w_j^2 \right] + \gamma T$$

Then for fixed structure $q(x)$, the optimal weights are

$$w_j^* = - \frac{G_j}{H_j + \lambda}$$

$$\Rightarrow \mathcal{L}^{(t)}(q) = - \frac{1}{2} \sum_{j=1}^{T} \frac{G_j^2}{H_j + \lambda} + \gamma T$$

This can be used as a scoring function for evaluating splits:

$$gain = \frac{1}{2} \left[ \frac{G_L^2}{H_L + \lambda} + \frac{G_R^2}{H_R + \lambda} - \frac{(G_L + G_R)^2}{H_L + H_R + \lambda} \right] - \gamma$$

         score of left     Score of      score if     complexity cost
          child           right child    no split      of split.

---

Exact greedy algorithm for split finding

$I$ = instance set of current node

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

iterate over features → for $k = 1$ to $m$:
     $G_L = 0$, $H_L = 0$
     for $j$ in sorted $(I, $ by $x_{jk})$:
         $G_L += g_j$,   $H_L += h_j$
         $G_R = G - G_L$,   $H_R = H - H_L$     Left to right search for best split value
         score $=$ max (score, gain)

output the split with maximum score

We can also regularise with shrinkage, scaling new terms with a small parameter $\eta$. We can use column/row subsampling which has the added benefit of speeding up computation.

## Alternative split finding algorithms

The exact greedy algorithm enumerates all possible splits and is thus computationally demanding. We can reduce the search space based on the weighted quantik sketch algorithm, which will propose a set of splits $S_K = \{S_{K1}, S_{K2}, S_{K3}, \cdots S_{KL}\}$ for feature $K$.

Sparsity-aware split finding can be implemented by adding a default direction to each node, which will be learnt from non-missing data.

## XGBoost system design

- To reduce the cost of sorting: before any learning is done, each column is sorted by feature value then stored in a block.
- The gradient statistics stay in place, referenced with pointers.
- Blocks $\Rightarrow$ parallelisation, and the column structure facilitates feature subsampling.
- Using the presorted blocks, finding quantiles is a linear scan.
- Too small a block size leads to inefficient parallelisation, while too large a block size means that not all of the gradient statistics will fit into the CPU cache. $2^{16}$ examples per block is a good compromise.
- Out-of-core computation, i.e disk reading/IO often takes a lot of computation time. This is improved by
-   — block compression by column
    — block sharding onto multiple disks.