

What Functions Does XGBoost Learn?

Dohyeong Ki & Adityanand Guntuboyina

Department of Statistics, University of California, Berkeley

XGBoost

Although XGBoost (eXtreme Gradient Boosting) has achieved remarkable empirical success, it has not been theoretically well-understood yet.

XGBoost fits a **finite sum of regression trees** to data.

XGBoost aims to (approximately) minimize least squares plus

$$\gamma \sum_k T_k + \alpha \sum_k \|w_k\|_1$$

squared L^2 norm is also common

where (1) T_k is the number of leaves in the k th regression tree,
(2) w_k is its vector of leaf weights.

XGBoost produces a discrete-valued tree fit, but it seems it also learns continuous functions quite effectively.

Q. What kinds of functions does XGBoost learn well?

Function Class Extending Finite Sums of Trees

Every regression tree can be expressed as a **finite linear combination** of

$$b_{\mathbf{p},\mathbf{q},\mathbf{t}}^S(x_1, \dots, x_d) := \prod_{j \in S} \{ \mathbf{1}(q_j = 0) \mathbf{1}(p_j(x_j - t_j) \geq 0) + \mathbf{1}(q_j = 1) \mathbf{1}(p_j(x_j - t_j) > 0) \}$$

where (1) $S \subseteq \{1, \dots, d\}$, (2) each $q_j \in \{0, 1\}$, and (3) $p_j \in \{-1, 1\}$.

- q_j determines whether the inequality is **weak (\geq) or strict ($>$)**.
- p_j controls the **direction** of the inequality.
- t_j is a **threshold** associated with variable x_j .

Example: $d = 2$ and $S = \{1, 2\}$

(1) $(p_1, q_1) = (1, 0)$ and $(p_2, q_2) = (-1, 0)$,

$$b_{\mathbf{p},\mathbf{q},\mathbf{t}}^S(x_1, x_2) = \mathbf{1}(x_1 \geq t_1) \cdot \mathbf{1}(x_2 \leq t_2)$$

(2) $(p_1, q_1) = (1, 1)$ and $(p_2, q_2) = (-1, 1)$,

$$b_{\mathbf{p},\mathbf{q},\mathbf{t}}^S(x_1, x_2) = \mathbf{1}(x_1 > t_1) \cdot \mathbf{1}(x_2 < t_2)$$

We consider **infinite linear combinations** of these basis functions with $|S| \leq s$.

We define $\mathcal{F}_{\infty\text{-st}}^{d,s}$ as the collection of all functions $f: \mathbb{R}^d \rightarrow \mathbb{R}$ of the form:

$$f_{c, \{\nu_{\mathbf{p},\mathbf{q}}^S\}}(x_1, \dots, x_d) = c + \sum_{S: 0 < |S| \leq s} \sum_{\mathbf{p} \in \{-1, 1\}^{|S|}} \sum_{\mathbf{q} \in \{0, 1\}^{|S|}} \int_{\mathbb{R}^{|S|}} b_{\mathbf{p},\mathbf{q},\mathbf{t}}^S(x_1, \dots, x_d) d\nu_{\mathbf{p},\mathbf{q}}^S(t_j, j \in S)$$

$\mathcal{F}_{\infty\text{-st}}^{d,s}$ is an **infinite dimensional extension** of the class $\mathcal{F}_{\text{st}}^{d,s}$ of **finite sums of regression trees with maximum depth s** .

→ consistent with XGBoost whose **max_depth** = 6 by default.

Complexity Extending XGBoost Penalty

Infinite linear combination representation for $f \in \mathcal{F}_{\infty\text{-st}}^{d,s}$ is not unique.

Define the **complexity** of $f \in \mathcal{F}_{\infty\text{-st}}^{d,s}$ as

$$V_{\infty\text{-xgb}}^1(f) := \inf \left\{ \sum_{S: 0 < |S| \leq s} \sum_{\mathbf{p} \in \{-1, 1\}^{|S|}} \sum_{\mathbf{q} \in \{0, 1\}^{|S|}} \|\nu_{\mathbf{p},\mathbf{q}}^S\|_{\text{TV}} : f_{c, \{\nu_{\mathbf{p},\mathbf{q}}^S\}} \equiv f \right\}$$

where $\|\nu\|_{\text{TV}}$ denotes the total variation of a signed measure ν .

Main Result 1:

If $f \in \mathcal{F}_{\text{st}}^{d,s}$, i.e., f is a **finite sum of regression trees**,

$$V_{\infty\text{-xgb}}^1(f) = \inf \left\{ \sum_k \|w_k\|_1 \right\}$$

where the infimum is over all representations of f in a finite sum of trees.

→ $V_{\infty\text{-xgb}}^1(\cdot)$ is an **extension** of the XGBoost penalty **with $\gamma = 0$**

$\gamma = 0$ means **no penalty on numbers of leaves**; the default choice by XGBoost

Relation to Hardy–Krause Variation

As the domain \mathbb{R}^d is unbounded, we need to **place an anchor** for Hardy–Krause variation **at infinity** (either $-\infty$ or $+\infty$ for each coordinate).

Let $\mathbf{a} = (a_1, \dots, a_d) \in \{-\infty, \infty\}^d$ denote the **anchoring point**.

For a function $f: \mathbb{R}^d \rightarrow \mathbb{R}$ and $S \subseteq \{1, \dots, d\}$, define

$$f_{(a_j, j \in S^c)}^S(x_j, j \in S) = \lim_{(x_j, j \in S^c) \rightarrow (a_j, j \in S^c)} f(x_1, \dots, x_d) \text{ for } (x_j, j \in S) \in \mathbb{R}^{|S|}$$

Hardy–Krause variation of f anchored at \mathbf{a} is defined by

$$\text{HK}_{\mathbf{a}}(f) = \sum_{\emptyset \neq S \subseteq \{1, \dots, d\}} \text{Vit}(f_{(a_j, j \in S^c)}^S).$$

where $\text{Vit}(\cdot)$ denotes Vitali variation.

Hardy–Krause variation is **asymmetric**, whereas $V_{\infty\text{-xgb}}^1(\cdot)$ is **symmetric**.

Example: $d = s = 2$ and $\mathbf{a} = (-\infty, -\infty)$

$$\text{HK}_{\mathbf{a}}(\mathbf{1}(\cdot_1 \geq t_1, \cdot_2 \geq t_2)) = 1 \text{ but } \text{HK}_{\mathbf{a}}(\mathbf{1}(\cdot_1 < t_1, \cdot_2 < t_2)) = 3$$

$$V_{\infty\text{-xgb}}^1(\mathbf{1}(\cdot_1 \geq t_1, \cdot_2 \geq t_2)) = V_{\infty\text{-xgb}}^1(\mathbf{1}(\cdot_1 < t_1, \cdot_2 < t_2)) = 1$$

In fact, $V_{\infty\text{-xgb}}^1(\cdot)$ is a **symmetrized version** of Hardy–Krause variation;

$V_{\infty\text{-xgb}}^1(\cdot)$ is the **infimal convolution** of $\text{HK}_{\mathbf{a}}(\cdot)$ over all anchors $\mathbf{a} \in \{-\infty, \infty\}^d$, when restricted to the subclass $\mathcal{F}_{\infty\text{-rst}}^{d,s}$ consisting of right-continuous functions

$$V_{\infty\text{-xgb}}^1(f) = \inf \left\{ \sum_{\mathbf{a} \in \{-\infty, \infty\}^d} \text{HK}_{\mathbf{a}}(f_{\mathbf{a}}) : \sum_{\mathbf{a} \in \{-\infty, \infty\}^d} f_{\mathbf{a}} \equiv f \text{ and } f_{\mathbf{a}} \in \mathcal{F}_{\infty\text{-rst}}^{d,s} \right\}$$

Least Squares Estimator (LSE)

A central object of interest is the following **least squares estimator**:

$$\underset{f}{\operatorname{argmin}} \left\{ \sum_{i=1}^n (y_i - f(\mathbf{x}^{(i)}))^2 : f \in \mathcal{F}_{\infty\text{-st}}^{d,s} \text{ and } V_{\infty\text{-xgb}}^1(f) \leq V \right\}.$$

Let $\mathcal{F}_{\text{rstm}}^{d,s}$ denote the collection of all **finite linear combinations** of $b_{\mathbf{p},\mathbf{q},\mathbf{t}}^S$ where

• $|S| \leq s \rightarrow$ depth is no larger than s

• $\mathbf{p} = \mathbf{1} - 2\mathbf{q} \rightarrow b_{\mathbf{p},\mathbf{q},\mathbf{t}}^S$ are **products only of $\mathbf{1}(x_j \geq t_j)$ and $\mathbf{1}(x_j < t_j)$**

→ aligns with XGBoost's tree-splitting scheme where one branch corresponds to $\mathbf{1}(x_j \geq t_j)$ and the other to $\mathbf{1}(x_j < t_j)$

• t_j are **midpoints between observed values** of the j th covariate

→ aligns with XGBoost's split points for numerical variables

→ By default (**tree_method** = **auto**), XGBoost uses midpoints when datasets are small but switches to quantiles for larger datasets.

Main Result 2:

The least squares estimator $\hat{f}_{n,V}^{d,s}$ over all $f \in \mathcal{F}_{\text{rstm}}^{d,s}$ with $V_{\infty\text{-xgb}}^1(f) \leq V$ is a least squares estimator over all $f \in \mathcal{F}_{\infty\text{-st}}^{d,s}$ with $V_{\infty\text{-xgb}}^1(f) \leq V$.

XGBoost can be viewed as a **greedy solver** for the **penalized version** of this **least squares problem over $\mathcal{F}_{\text{rstm}}^{d,s}$** .

Theoretical Accuracy of LSE

Main Result 3:

Assume the following **random design** setting:

(1) $y_i = f^*(\mathbf{x}^{(i)}) + \epsilon_i$ where $f^* \in \mathcal{F}_{\infty\text{-st}}^{d,s}$ and $\epsilon_i \stackrel{\text{i.i.d.}}{\sim} N(0, \sigma^2)$ can be replaced by a more general assumption

(2) $\mathbf{x}^{(i)} \stackrel{\text{i.i.d.}}{\sim} p_0$ for some density p_0

(3) p_0 has compact support: there exist $M_1, \dots, M_d > 0$ such that

$$p_0(\mathbf{x}) = 0 \text{ unless } \mathbf{x} \in \prod_{j=1}^d \left[-\frac{M_j}{2}, \frac{M_j}{2} \right]$$

(4) p_0 is bounded above; $B := M_1 \cdots M_d \cdot \sup_{\mathbf{x}} p_0(\mathbf{x}) < +\infty$.

If $V > V_{\infty\text{-xgb}}^1(f^*)$, then we have

$$\mathbb{E} \left[\int (\hat{f}_{n,V}^{d,s}(\mathbf{x}) - f^*(\mathbf{x}))^2 \cdot p_0(\mathbf{x}) d\mathbf{x} \right] = O(n^{-2/3} (\log n)^{4(s-1/3)}).$$

constant factor depends on B, d, V , and σ

It can also be proved that this rate is **nearly minimax optimal**.

Whether XGBoost itself achieves a similar nearly dimension-free rate of convergence is an open problem.