# XBART ACCELERATED BAYESIAN ADDITIVE REGRESSION TREES

## FASTER and MORE ACCURATE than XGBoost

R & python code now available. Coming soon on CRAN and pip.

BART boasts state-of-the-art prediction accuracy. But, **BART MCMC can be SLOW.**

$$\frac{1}{2} \sum_{b=1}^{B} \left\{ \log \left( \frac{\sigma^2}{\sigma^2 + \tau n_b} \right) + \frac{\tau}{\sigma^2 (\sigma^2 + \tau n_b)} s_b^2 \right\}$$

**XBART** grows trees stochastically but recursively, using the unique BART split criteria, so it **is FAST.**

By growing trees recursively, many efficiency tricks can be exploited: pre-sorting variables, adaptive nested cutpoints, sparse trees.
Scan the QR code to see the paper for details.

by Jingyu He, Saar Yalov and P. Richard Hahn

# XBART: Accelerated Bayesian Additive Regression Trees

Jingyu He [1]    Saar Yalov [2]    P. Richard Hahn [2]

[1]University of Chicago, jingyuhe@chicagobooth.edu    [2]Arizona State University

## Highlights

XBART is motivated by Bayesian additive regression trees (BART), provides fast posterior estimation for BART model. Simulation shows that

1. XBART is faster and more accurate than `xgboost` with tuning parameters by cross validation.
2. Fit large data set (250K observations) in tolerable time, which BART can never do.

## BART Prior

Bayesian Additive Regression Trees, first appeared in Chipman et al. (2010). BART is not merely a version of random forest or boosted regression trees in which prior distributions have been placed over model parameters, but **prior over tree structure** and **parameters**.

Pros Robust to tuning parameter, **more accurate prediction**, a natural Bayesian measure of uncertainty.

Cons The random walk Metropolis-Hastings Markov chain Monte Carlo algorithm is slow.

The BART model is

$$y = \sum_{l=1}^{L} g_l(x, T_l, \mu_l) + \epsilon \qquad (1)$$

where $T_l$ denotes regression tree and $\mu_l$ is vector of means associated to all nodes of tree $l$. The BART prior has three components

1. Probability of a node having children at depth $d$ is
$$\alpha(1+d)^{-\beta}$$

2. Uniform distribution over available predictors to split at.

3. Uniform distribution on a discrete set of available splitting values for the assigned predictor.

The basic BART MCMC takes a Metropolis-within-Gibbs algorithm, update each tree by local random walk Metropolis-Hastings (MH) update. **Slow, cannot work on large data set.**

## Bayesian Backfitting

Taking advantage of the additive structure of the model, these updates can be written as

1. $T_l, \mu_l \mid r_l, \sigma^2$, for $l = 1, \ldots, L$, which is done compositionally (for each $l$) as
   1. $T_l \mid r_l, \sigma^2$,
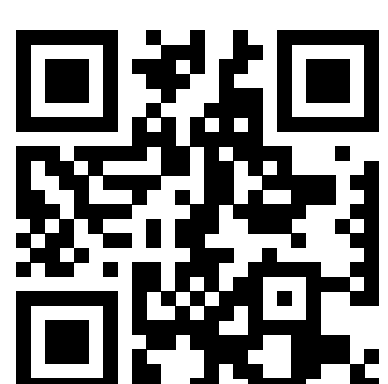   2. $\mu_l \mid T_l, r_l, \sigma^2$,
2. $\sigma^2 \mid r$.

for "residuals" defined as

$$r_l^{(k+1)} \equiv y - \sum_{l' < l} g(\mathbf{X}; T_{l'}, \mu_{l'})^{(k+1)} - \sum_{l' > l} g(\mathbf{X}; T_{l'}, \mu_{l'})^{(k)},$$

and

$$r^{(k)} \equiv y - \sum_{l=1}^{L} g(\mathbf{X}; T_l, \mu_l)^{(k)},$$

where $k$ indexes the Monte Carlo iteration.

Code available upon request for now, will be on CRAN and pip (python version) soon.

## Grow-from-root Backfitting

Given the current node, the likelihood of each cut-point candidate is

$$\pi(v, c) = \frac{\exp\left(\ell(c, v)\right)\kappa(c)}{\sum_{v'=1}^{V}\sum_{c'=0}^{C}\exp\left(\ell(c', v')\right)\kappa(c')} \qquad (2)$$

where

$$\ell(v,c) = \frac{1}{2}\left\{\log\left(\frac{\sigma^2}{\sigma^2 + \tau n(\leq, v, c)}\right) + \frac{\tau}{\sigma^2(\sigma^2 + \tau n(\leq, v, c))}s(\leq, v, c)^2\right\}$$
$$+ \frac{1}{2}\left\{\log\left(\frac{\sigma^2}{\sigma^2 + \tau n(>, v, c)}\right) + \frac{\tau}{\sigma^2(\sigma^2 + \tau n(>, v, c))}s(>, v, c)^2\right\}$$

for $c \neq 0$. $n(\leq, v, c)$ is the number of observations in the current leaf node that have $x_v \leq c$ and $s(\leq, v, c)$ is the sum of the residual $r_l^{(k)}$; $n(>, v, c)$ and $s(>, v, c)$ are defined analogously. Also, $\kappa(c \neq 0) = 1$.

For $c = 0$, corresponding to **stop-splitting option**, we have instead

$$\ell(v,c) = \frac{1}{2}\left\{\log\left(\frac{\sigma^2}{\sigma^2 + \tau n}\right) + \frac{\tau}{\sigma^2(\sigma^2 + \tau n)}s^2\right\}$$

and $\kappa(0) = \frac{1 - \alpha(1+d)^{-\beta}}{\alpha(1+d)^{-\beta}}$, where $n = n(\leq, v, c) + n(>, v, c)$, $s = s(\leq, v, c) + s(<, v, c)$.

## Pre-sorting Features for Efficiency

Observe that the BART criterion **depends on the partition sum only**. With sorted predictor variables, the likelihood of cut-point can be computed via a single sweep through the data (per variable), taking cumulative sum.

$$s(\leq, v, c) = \sum_{h \leq c} r_{o_{vh}}$$

and

$$s(>, v, c) = \sum_{h=1}^{n} r_{lh} - s(\leq, v, c).$$

## Recursive cut-points

Take every $j$th value (starting from the smallest) as an eligible split point with $j = \lfloor \frac{n_b - 2}{C} \rfloor$.

As a default, we set the number of cut-points to $\max(\sqrt{n}, 100)$, where $n$ is the sample size of the entire data set.

## Sparse Trees

We considering $m \leq V$ variables at a time when sampling each splitting rule.

We introduce a parameter vector w which denotes the prior probability that a given variable is chosen to be split on, as suggested in Linero (2016).

**Before sampling each splitting rule, we randomly select $m$ variables with probability proportional to** w. These $m$ variables are sampled sequentially and *without replacement*, with selection probability proportional to w.

## Variable Importance Weights

The variable weight parameter w is given a Dirichlet prior with hyperparameter $\bar{w}$ set to all ones.

Split counts are updated in between each tree sampling/growth step:

$$\bar{w} \leftarrow \bar{w} - \bar{w}_l^{(k-1)} + \bar{w}_l^{(k)} \qquad (3)$$

where $\bar{w}_l^{(k)}$ denotes the length-$V$ vector recording the number of splits on each variable in tree $l$ at iteration $k$. The weight parameter is then resampled as w $\sim$ Dirichlet($\bar{w}$).

## Posterior Prediction

Given $K$ iterations of the algorithm, suppose $I < K$ is denotes the length of the burn-in period, the final prediction is

$$\bar{f}(\mathbf{X}) = \frac{1}{K - I}\sum_{k > I}^{K} f^{(k)}(\mathbf{X}). \qquad (4)$$

where $f^{(k)}$ denotes a sample of the forest.

## XBART algorithm

**Algorithm 1** Grow-from-root backfitting

1: $N \leftarrow$ number of rows of $y, x$
2: Sample $m$ variables use weight $w$ as shown in section sparse trees.
3: Select $C$ cutpoints as shown in section Grow-from-root backfitting.
4: Evaluate $C \times m + 1$ candidate cutpoints and no-split option with equation (2).
5: Sample one cutpoint propotional to equation (2).
6: **if** sample no-split option **then**
7:     Sample leaf parameter from normal distribution $\mu \sim N\left(\sum y / \left[\sigma^2\left(\frac{1}{\tau} + \frac{N}{\sigma^2}\right)\right], 1/\left[\frac{1}{\tau} + \frac{N}{\sigma^2}\right]\right)$. **return**
8: **else**
9:     $w_l[j] = w_l[j] + 1$, add count of selected split variable.
10:     Split data to left and right node.
11:     GROW_FROM_ROOT($y_{\text{left}}, \mathbf{X}_{\text{left}}, C, m, \text{w}, \sigma^2$)
12:     GROW_FROM_ROOT($y_{\text{right}}, \mathbf{X}_{\text{right}}, C, m, \text{w}, \sigma^2$)
13: **end if**

**Algorithm 2** Accelerated Bayesian Additive Regression Trees (XBART)

   $V \leftarrow$ number of columns of $\mathbf{X}$
2: $N \leftarrow$ number of rows of $\mathbf{X}$
   Initialize $r_l^{(0)} \leftarrow y/L$.
4: **for** $k$ in 1 to $K$ **do**
   **for** $l$ in 1 to $L$ **do**
6:     Calculate residual $r_l^{(k)}$ as shown in section Bayesian Backfitting.
       **if** $k < I$ **then**
8:         GROW_FROM_ROOT($r_l^{(k)}, \mathbf{X}, C, V, \text{w}, \sigma^2$) {use all variables in burnin iterations}
       **else**
10:         GROW_FROM_ROOT($r_l^{(k)}, \mathbf{X}, C, m, \text{w}, \sigma^2$)
       **end if**
12:     $\bar{w} \leftarrow \bar{w} - \bar{w}_l^{(k-1)} + \bar{w}_l^k$ {update $\bar{w}$ with split counts of current tree}
       w $\sim$ Dirichlet($\bar{w}$)
14:     $\sigma^2 \sim$ Inverse-Gamma($N + \alpha, r_l^{(k)t}r_l^{(k)} + \eta$)
   **end for**
16: **end for**
   **return**

## Is it a valid MCMC algorithm?

The algorithm works well on its own right. We can use it **as proposal of M-H algorithm**, rather than a random walk M-H, to get full Bayesian inference. Future work.

## Simulations

| $n$ | XBART | XGB+CV | XGB | NN |
|---|---|---|---|---|
| | | Linear | | |
| 10k | 5.07 (16) | 8.04 (61) | 21.25 (0) | 7.39 (12) |
| 50k | 3.16 (135) | 5.47 (140) | 16.17 (4) | 3.62 (14) |
| 250k | 2.03 (1228) | 3.15 (1473) | 11.49 (54) | 1.89 (19) |
| | | Max | | |
| 10k | 1.94 (16) | 2.76 (60) | 7.18 (0) | 2.98 (15) |
| 50k | 1.22 (133) | 1.85 (139) | 5.49 (4) | 1.63 (16) |
| 250k | 0.75 (1196) | 1.05 (1485) | 3.85 (54) | 0.85 (22) |
| | | Single Index | | |
| 10k | 7.13 (16) | 10.61 (61) | 28.68 (0) | 9.43 (14) |
| 50k | 4.51 (133) | 6.91 (139) | 21.18 (4) | 6.42 (16) |
| 250k | 3.06 (1214) | 4.10 (1547) | 14.82 (54) | 4.72 (21) |
| | | Trig + Poly | | |
| 10k | 4.94 (16) | 7.16 (61) | 17.97 (0) | 8.20 (13) |
| 50k | 3.01 (132) | 4.92 (139) | 13.30 (4) | 5.53 (14) |
| 250k | 1.87 (1216) | 3.17 (1462) | 9.37 (49) | 4.13 (20) |

Table 1. Root mean squared error (running time).