

Pushing the Envelope of Gradient Boosting Forests via Globally-Optimized Oblique Trees

Magzhan Gabidolla and Miguel Á. Carreira-Perpiñán, Dept. Computer Science & Engineering, UC Merced

L Abstract

Ensemble methods based on decision trees, such as Random Forests or boosted class 1 forests, have long been established as some of the most powerful, off-the-shelf machine learning models, and have been widely used in computer vision and other areas. In recent years, a specific form of boosting, gradient boosting (GB), \bigcirc has gained prominence. This is partly because of highly optimized implementations such as XGBoost or LightGBM, which incorporate many clever modifications and heuristics. However, one gaping hole remains unexplored in GB: the $\Delta = 10$ $\Delta = 40'$ Щ 10 construction of individual trees. To date, all successful GB versions use axis-TAO $\Delta = 2$ aligned trees trained in a suboptimal way via greedy recursive partitioning. We TAO $\Delta = 2$ address this gap by using a more powerful type of trees (having hyperplane splits) 200 Number of trees T in SPORF Number of trees *T* in XGBoost and an algorithm that can optimize, globally over all the tree parameters, the objective function that GB dictates. We show, in several benchmarks of image and Consider a synthetic binary classification of MNIST digit images where class 1 satisfies that other data types, that GB forests of these stronger, well-optimized trees consis- $(Q_1 + Q_4) - (Q_2 + Q_3) \ge P$ or $(Q_2 + Q_3) - (Q_1 + Q_4) \ge P$, where the 28×28 pixel image is tently exceed the test accuracy of axis-aligned forests from XGBoost, LightGBM split into 4 quadrants $\begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}$ and Q_i is the sum of the [0, 1] pixel intensities in quadrant i, and and other strong baselines. Further, this happens using many fewer trees and P = 30. Class 1 is an "(anti)diagonally dominant" image like \blacksquare or \blacksquare . A single depth-2 tree sometimes even fewer parameters overall. trained with TAO achieves a near-perfect training/test error of 0.0%/0.96%. XGBoost forests Work supported by NSF award IIS–2007147 do much worse. A recent algorithm to train oblique forests (SPORF) performs far worse than **2** Gradient Boosting (GB) Forests the XGBoost, suggesting that proper oblique tree optimization is essential.

- Ensembles of decision trees have long been established as some of the most powerful, off-the-shelf machine learning models.
- In recent years, one type of forest, Gradient Boosting (GB), has gained prominence due to their:
- Strong empirical performance on many problems
- The development of extremely efficient implementations such as XGBoost or LightGBM.
- They typically require little effort on hyperparameter tuning and are thus con sidered "off-the-shelf".
- Given the tremendous effort put on the development and refinement of the popular GB toolkits, how can we further improve GB forests?

Nodeling high-order feature interactions: O axis-aligned trees oblique trees VS $x^{T}w_{1} < 3.2$ *x*₃₂ < 1.2 *x*₅₁ < -2.0 $\langle {\bf x}^{T} {\bf w}_{2} \geq 0.5$ $x_2 \ge 3.4$ Each decision node is a function of all

Only 3 features participate in the routing function of the above leaf. Maximum order of feature interactions is limited by the depth Δ in axis-aligned trees.

the features. Their non-linear combination is a much complex order-D interaction. For modeling complex functions, a forest of oblique trees should achieve higher accuracy and require tewer trees.



4 Learning oblique trees in the GB framework

To learn oblique trees in GB, we build on a recent algorithm, tree alternating optimization (TAO). Given a decision tree of some predetermined structure $\tau(\mathbf{x}; \Theta)$ with parameters $\Theta = \{(\mathbf{w}_i, w_{i0})\}_{i \in D} \cup \{\theta_i\}_{i \in L}, \text{ decision nodes in set } D, \text{ leaves in set } L, \text{ we optimize:} \}$

TAO is based on two theorems. First, eq. (1) separates over any subset of non-descendant nodes (e.g. all the nodes at the same depth); this follows from the fact that the tree makes hard decisions. Second, optimizing over the parameters of a single node *i* simplifies to a well-defined reduced problem over the instances that currently reach node i (the reduced set $\mathcal{R}_i \subset \{1, \ldots, N\}$). The form of the reduced problem depends on the type of node: Decision node It is a weighted 0/1 loss binary classification problem, where the two classes correspond to the left and right child, which are the only possible outcomes for an instance. Child left; (right;) incurs a loss (weight) given by the prediction of the leaf reached from the left (right) child's subtree. Thus, each instance is assigned as pseudolabel the child with lower loss. This problem is NP-hard but can be well approximated with a convex surrogate; we use ℓ_1 -regularized logistic regression where each instance is weighted by the loss difference between the winner and the other child. Leaf The reduced problem consists of optimizing the original loss but over the leaf classifier

on its reduced set:

$$\min_{\boldsymbol{\theta}_i} \sum_{n \in \mathcal{R}_i} \mathbf{g}_n^T \boldsymbol{\theta}_i + \frac{1}{2} \boldsymbol{\theta}_i^T \mathbf{H}_n \boldsymbol{\theta}_i.$$

This is solved similarly as in XGBoost and other frameworks.

Given an initial tree structure with initial parameter values, the resulting algorithm repeatedly visits nodes in reverse breadth-first search order. Each iteration trains all nodes at the same depth (in parallel) from the leaves to the root, by solving either an ℓ_1 -regularized logistic regression at each decision node, or the above exact solution as each leaf.



$$\min \sum I(\mathbf{g}_n, \mathbf{H}_n, \boldsymbol{\tau}(\mathbf{x}_n; \boldsymbol{\Theta})) + \alpha \sum \|\mathbf{w}_i\|_1 \quad \text{with} \quad I(\mathbf{g}, \mathbf{H}, \boldsymbol{\gamma}) = \mathbf{g}^T \boldsymbol{\gamma} + \frac{1}{2} \boldsymbol{\gamma}^T \mathbf{H} \boldsymbol{\gamma}.$$
(1)

Pendigits					News20					MNIST				
Etest (%)	#pars.	Т	Δ	Forest	Etest (%)	#pars.	Т	Δ	Forest	Etest (%)	#pars.	T	Δ	
$3.49{\pm}0.00$	90k	1k	11	GB-sklearn	$23.42{\pm}0.03$	156k	2k	6	XGBoost	$4.38{\pm}0.00$	70k	100	10	
$3.46{\pm}0.00$	18k	1k	4	SPORF	$22.51{\pm}0.09$	(1.3M)	100	569	GB-TAO	$4.17{\pm}0.08$	21k	1	12	
$3.46{\pm}0.00$	137k	10k	4	XGBoost	$21.39{\pm}0.00$	705k	20k	6	SPORF	$2.89{\pm}0.04$	(143M)	1k	50	
$3.31 {\pm} 0.00$	895k	10k	4	XGBoost	$21.34{\pm}0.00$	188k	6k	6	XGBoost	$2.20{\pm}0.00$	107k	1k	6	
$3.15{\pm}0.25$	1.3k	1	8	LightGBM	$20.69 {\pm} 0.00$	1.8M	20k	27	LightGBM	$2.02{\pm}0.00$	121k	1k	10	
$2.91{\pm}0.09$	(1.6M)	1k	20	LightGBM	$19.78 {\pm} 0.00$	546k	6k	28	GB-TAO	$1.94{\pm}0.00$	671k	30	10	
$2.87{\pm}0.01$	(105k)	100	20	GB-TAO	$18.13{\pm}0.01$	479k	400	4	XGBoost	$1.91 {\pm} 0.00$	505k	10k	6	
$2.17{\pm}0.02$	13k	10	7	GB-TAO	$18.76 {\pm} 0.01$	746k	50	6	LightGBM	$1.62{\pm}0.00$	642k	10k	21	
$2.00{\pm}0.04$	44k	30	7	GB-TAO	$16.65 {\pm} 0.04$	1.6M	800	4	GB-TAO	$1.55{\pm}0.02$	7.2M	1.4k	7	
CT-slice				casp					Year					
Etest (%)	#pars.	Т	Δ	Forest	Etest (%)	#pars.	Т	Δ	Forest	E _{test}	#pars.	T	Δ	
$1.53{\pm}0.00$	153k	100	97	XGBoost	$3.66{\pm}0.00$	119k	100	10	GB-TAO	9.17±0.01	19k	1	8	
$1.52{\pm}0.00$	91k	1k	23	XGBoost	$3.58{\pm}0.00$	793k	1k	10	XGBoost	$9.05{\pm}0.00$) 153k	100	10	
$1.50 {\pm} 0.00$	107k	100	10	GB-sklearn	$3.58{\pm}0.01$	854k	1k	10	LightGBM	$9.03{\pm}0.00$) 153k	100	37	
$1.28{\pm}0.02$	28k	1	8	LightGBM	$3.54{\pm}0.00$	153k	100	114	GB -sklearr	ם 8.96±0.02	2 171k	1k	6	
$1.26{\pm}0.03$	900k	1k	10	GB-TAO	$3.49{\pm}0.01$	256k	50	12	LightGBM	$8.92{\pm}0.00$	1.5M	1k	43	
$1.26{\pm}0.00$	767k	1k	10	LightGBM	$3.48{\pm}0.00$	766k	1k	109	XGBoost	$8.91{\pm}0.00$	1.8M	1k	10	
$0.90{\pm}0.02$	81k	30	4	GB-TAO	$3.43{\pm}0.00$	481k	100	12	GB-TAO	8.88±0.02	2 78k	20	6	
$0.45{\pm}0.01$	1.2M	100	6	GB-TAO	$3.39{\pm}0.01$	887k	200	12	GB-TAO	8.73±0.01	402k	100	6	
	Pendig E_{test} (%) 3.49 ± 0.00 3.46 ± 0.00 3.46 ± 0.00 3.46 ± 0.00 3.15 ± 0.25 2.91 ± 0.09 2.87 ± 0.01 2.17 ± 0.02 2.00 ± 0.04 CT-slic E_{test} (%) 1.53 ± 0.00 1.52 ± 0.00 1.50 ± 0.00 1.26 ± 0.03 1.26 ± 0.03 0.45 ± 0.01	Pendigits E_{test} (%)#pars. 3.49 ± 0.00 90k 3.46 ± 0.00 137k 3.46 ± 0.00 137k 3.31 ± 0.00 895k 3.15 ± 0.25 1.3k 2.91 ± 0.09 (1.6M) 2.87 ± 0.01 (105k) 2.17 ± 0.02 13k 2.00 ± 0.04 44kCT-slice E_{test} (%)#pars. 1.53 ± 0.00 153k 1.52 ± 0.00 91k 1.50 ± 0.00 107k 1.28 ± 0.02 28k 1.26 ± 0.03 900k 1.26 ± 0.00 767k 0.90 ± 0.02 81k 0.45 ± 0.01 1.2M	Pendigits E_{test} (%)#pars. T 3.49 ± 0.00 90k1k 3.46 ± 0.00 18k1k 3.46 ± 0.00 137k10k 3.46 ± 0.00 137k10k 3.15 ± 0.25 1.3k1 2.91 ± 0.09 (1.6M)1k 2.87 ± 0.01 (105k)100 2.17 ± 0.02 13k10 2.00 ± 0.04 44k30CT-slice E_{test} (%)#pars. T 1.53 ± 0.00 153k100 1.52 ± 0.00 91k1k 1.50 ± 0.00 107k100 1.28 ± 0.02 28k1 1.26 ± 0.03 900k1k 0.90 ± 0.02 81k30 0.45 ± 0.01 1.2M100	Pendigits E_{test} (%)#pars. T Δ 3.49 ± 0.00 90k1k11 3.46 ± 0.00 18k1k4 3.46 ± 0.00 137k10k4 3.46 ± 0.00 137k10k4 3.31 ± 0.00 895k10k4 3.15 ± 0.25 1.3k18 2.91 ± 0.09 (1.6M)1k20 2.87 ± 0.01 (105k)10020 2.17 ± 0.02 13k107 2.00 ± 0.04 44k307CT-slice E_{test} (%)#pars. T Δ 1.53 ± 0.00 153k10097 1.52 ± 0.00 91k1k23 1.50 ± 0.00 107k10010 1.28 ± 0.02 28k18 1.26 ± 0.03 900k1k10 0.90 ± 0.02 81k304 0.45 ± 0.01 1.2M1006	PendigitsForest E_{test} (%)#pars. T Δ 3.49 ± 0.00 90k1k11 3.46 ± 0.00 18k1k4 3.46 ± 0.00 137k10k4 3.46 ± 0.00 137k10k4 3.46 ± 0.00 137k10k4 3.31 ± 0.00 895k10k4 3.31 ± 0.00 895k10k4 2.91 ± 0.09 (1.6M)1k20 2.87 ± 0.01 (105k)10020 2.87 ± 0.01 (105k)10020 2.00 ± 0.04 44k307 2.00 ± 0.04 44k307 E_{test} (%)#pars. T Δ E_{test} (%)#pars. T Δ 1.53 ± 0.00 153k10097 1.52 ± 0.00 91k1k23 1.50 ± 0.00 107k10010 1.28 ± 0.02 28k18 1.26 ± 0.03 900k1k10 0.90 ± 0.02 81k304 0.45 ± 0.01 1.2M1006	News2 E_{test} (%)#pars. T Δ Forest E_{test} (%) 3.49 ± 0.00 90k1k11GB-sklearn 23.42 ± 0.03 3.46 ± 0.00 13k1k4SPORF 22.51 ± 0.09 3.46 ± 0.00 137k10k4XGBoost 21.39 ± 0.00 3.31 ± 0.00 895k10k4XGBoost 21.34 ± 0.00 3.15 ± 0.25 1.3k18LightGBM 20.69 ± 0.00 2.91 ± 0.09 (1.6M)1k20LightGBM 19.78 ± 0.00 2.87 ± 0.01 (105k)10020GB-TAO 18.13 ± 0.01 2.17 ± 0.02 13k107GB-TAO 18.76 ± 0.01 2.00 ± 0.04 44k307GB-TAO 16.65 ± 0.04 CT-slicecaspEtest (%)#pars. T Δ Forest E_{test} (%) 1.53 ± 0.00 153k10097XGBoost 3.66 ± 0.00 1.52 ± 0.00 91k1k23XGBoost 3.58 ± 0.01 1.28 ± 0.02 28k18LightGBM 3.54 ± 0.00 1.26 ± 0.03 900k1k10GB-TAO 3.49 ± 0.01 1.26 ± 0.00 767k1k10LightGBM 3.48 ± 0.00 0.90 ± 0.02 81k304GB-TAO 3.43 ± 0.00 0.90 ± 0.02 81k304GB-TAO 3.43 ± 0.00	$\begin{array}{c c c c c c c c c c c c c c c c c c c $	PendigitsNews20 E_{test} (%)#pars. T Δ Forest E_{test} (%)#pars. T 3.49 ± 0.00 90k1k11GB-sklearn 23.42 ± 0.03 156k2k 3.46 ± 0.00 18k1k4SPORF 22.51 ± 0.09 $(1.3M)$ 100 3.46 ± 0.00 137k10k4XGBoost 21.39 ± 0.00 705k20k 3.31 ± 0.00 895k10k4XGBoost 21.34 ± 0.00 188k6k 3.15 ± 0.25 1.3k18LightGBM 20.69 ± 0.00 1.8M20k 2.91 ± 0.09 $(1.6M)$ 1k20LightGBM 19.78 ± 0.00 546k6k 2.87 ± 0.01 $(105k)$ 10020GB-TAO 18.13 ± 0.01 479k400 2.17 ± 0.02 13k107GB-TAO 18.76 ± 0.01 746k50 2.00 ± 0.04 44k307GB-TAO 16.65 ± 0.04 1.6M800CT-slicecaspCaspLightGBM 3.58 ± 0.00 793k1k 1.53 ± 0.00 153k10097XGBoost 3.58 ± 0.00 793k1k 1.52 ± 0.00 91k1k23XGBoost 3.58 ± 0.01 854k1k 1.28 ± 0.02 28k18LightGBM 3.54 ± 0.00 153k100 1.26 ± 0.03 900k1k10GB-TAO 3.43 ± 0.00 766k1k 0.90	$\begin{array}{c c c c c c c c c c c c c c c c c c c $	$\begin{array}{c c c c c c c c c c c c c c c c c c c $	$\begin{array}{c c c c c c c c c c c c c c c c c c c $	$\begin{array}{c c c c c c c c c c c c c c c c c c c $	$\begin{array}{c c c c c c c c c c c c c c c c c c c $	

Table 1: Comparison of different forest-based models for classification (top) and regression (bottom), sorted by decreasing test error. We report 0-1 test error or RMSE E_{test} (mean±std over 5 repeats), and the number of parameters in the model. T refers to the number of trees. Δ is the max depth of the forest.



parallel processing on a shared memory system with 8 processors.



Figure 1: Comparison of different GB forests as a function of the number of boosting steps M (top) and time (bottom). All methods except GB-sklearn are trained using