

# Smaller, more accurate regression forests using tree alternating optimization

Arman Zharmagambetov and Miguel Á. Carreira-Perpiñán

Dept. of Computer Science and Engineering  
University of California, Merced

ICML, July 2020



- A forest (= ensemble of decision trees) is a powerful machine learning method that has been successfully applied in many applications: computer vision, speech processing, NLP, etc. They are often (part of) the winning methods in ML competitions and challenges.
- Some examples of forests:
  - *Random forests* train each tree independently on a different data sample (bagging).
  - *Boosted Trees* sequentially train trees on reweighted versions of the data.

In both cases, the forest prediction is obtained by weighted averaging or voting.

We focus on **regression forests**, where the prediction is a continuous scalar or vector, using **bagging**.

Random forests (and their variations) for regression have important advantages:

- high predictive accuracy
- few hyperparameters to set
- reasonably fast to train and can be scaled to large datasets
- considered to be robust against overfitting.

But they have an important disadvantage: **the individual trees they learn are far from accurate**. This is due to two reasons:

- Each tree is axis-aligned (a decision node tests for a single feature, e.g. “if  $x_7 \geq 3$  then go right”). This is a very restrictive model, particularly for particularly for correlated features.
- Standard tree learning algorithms, based on a greedy recursive tree growth (such as CART), are highly suboptimal.

There are exist few works that propose forests of more complex trees (Menze et al., 2000; Breiman, 2001; Frank & Kramer, 2004), but they mostly focus on classification rather than regression and improve marginally over conventional axis-aligned trees.

- We address both issues by:
  - using trees with more complex nodes (oblique, i.e., hyperplane)
  - using a better optimization algorithm to learn the tree.

The resulting forests are smaller, shallower and much more accurate, consistently over various datasets.

- We build on a recently proposed algorithm for learning classification trees, **Tree Alternating Optimization (TAO)** (Carreira-Perpiñán et al. 2018). TAO finds good approximate optima of an objective function over a tree with predetermined structure and it applies to trees beyond axis-aligned splits.
- We adapt TAO to the regression case and then use it in combination with bagging to learn forests of oblique trees.

We consider trees whose nodes make hard decisions (not soft trees). Optimizing such trees is difficult because they are not differentiable. Assuming a tree structure  $\mathbf{T}$  is given (say, binary complete of depth  $\Delta$ ), consider the following optimization problem over its parameters:

$$E(\Theta) = \sum_{n=1}^N L(\mathbf{y}_n, \mathbf{T}(\mathbf{x}_n; \Theta)) + \alpha \sum_{i \in \mathcal{N}} \phi_i(\theta_i)$$

given a training set  $\{(\mathbf{x}_n, \mathbf{y}_n)\}_{n=1}^N$ .  $\Theta = \{\theta_i\}_{i \in \mathcal{N}}$  is a set of parameters of all tree nodes. The loss function  $L(\mathbf{y}, \mathbf{z})$  is the squared error  $\|\mathbf{y} - \mathbf{z}\|_2^2$  in our case (although it is possible to use other losses, such as the least absolute deviation or a robust loss). The regularization term  $\phi_i$  (e.g.  $\ell_1$  norm) penalizes the parameters  $\theta_i$  of each node.

Our TAO algorithm for regression is based on 3 theorems: separability condition, reduced problem over a leaf, reduced problem over a decision node.

## 1. Separability condition

Consider any pair of nodes  $i$  and  $j$ . Assume the parameters of all other nodes ( $\Theta_{\text{rest}}$ ) are fixed. If nodes  $i$  and  $j$  are not descendants of each other, then  $E(\Theta)$  can be rewritten as:

$$E(\Theta) = E_i(\theta_i) + E_j(\theta_j) + E_{\text{rest}}(\Theta_{\text{rest}})$$

In other words, the separability condition states that any set of non-descendant nodes of a tree can be optimized independently. Note that  $E_{\text{rest}}(\Theta_{\text{rest}})$  can be treated as a constant since we fix  $\Theta_{\text{rest}}$ .

A set of non-descendant nodes are all the leaves. Optimizing over the parameters of one leaf is given by the following theorem.

## 2. Reduced problem over a leaf

If  $i$  is a leaf, the optimization of  $E(\Theta)$  over  $\theta_i$  can be equivalently written as:

$$\min_{\theta_i} E_i(\theta_i) = \sum_{n \in \mathcal{R}_i} L(\mathbf{y}_n, \mathbf{g}_i(\mathbf{x}_n; \theta_i)) + \alpha \phi_i(\theta_i)$$

The **reduced set**  $\mathcal{R}_i$  contains the training instances that reach leaf  $i$ . Each leaf  $i$  has a predictor function  $\mathbf{g}_i(\mathbf{x}; \theta_i): \mathbb{R}^D \rightarrow \mathbb{R}^K$  (we use a constant or linear regressor) that produces the actual output. Therefore, solving the reduced problem over a leaf  $i$  amounts to fitting the leaf's predictor  $\mathbf{g}_i$  to the instances in its reduced set to minimize the original loss (e.g. squared error).

An example of a set of non-descendant nodes are all the decision nodes at the same depth. Optimizing over the parameters of one decision node is given by the following theorem.

### 3. Reduced problem over a decision node

If  $i$  is a decision node, the optimization of  $E(\Theta)$  over  $\theta_i$  can be equivalently written as:

$$\min_{\theta_i} E_i(\theta_i) = \sum_{n \in \mathcal{R}_i} l_{in}(f_i(\mathbf{x}_n; \theta_i)) + \alpha \phi_i(\theta_i)$$

where  $\mathcal{R}_i$  is the reduced set of node  $i$  and (assuming binary trees)  $f_i(\mathbf{x}; \theta_i): \mathbb{R}^D \rightarrow \{\text{left}, \text{right}\}$  is a decision function in node  $i$  which sends instance  $\mathbf{x}_n$  to the corresponding child of  $i$ . We consider oblique trees, having hyperplane decision functions “go to right if  $\mathbf{w}_i^T \mathbf{x} + w_{i0} \geq 0$ ” (where  $\theta_i = \{\mathbf{w}_i, w_{i0}\}$ ).  $l_{in}(\cdot)$  is the loss incurred if  $\mathbf{x}_n$  chooses the right or left subtree.

# Optimizing a single tree with TAO: decision nodes (cont.)

The reduced problem over a decision node can be equivalently rewritten as a weighted 0/1 loss binary classification problem on the node's reduced set instances:

$$\min_{\theta_i} E_i(\theta_i) = \sum_{n \in \mathcal{R}_i} \bar{L}_{in}(\bar{y}_{in}, f_i(\mathbf{x}_n; \theta_i)) + \alpha \phi_i(\theta_i)$$

where the weighted 0/1 loss  $\bar{L}_{in}(\bar{y}_{in}, \cdot)$  for instance  $n \in \mathcal{R}_i$  is defined as  $\bar{L}_{in}(\bar{y}_{in}, y) = l_{in}(y) - l_{in}(\bar{y}_{in}) \forall y \in \{\text{left}, \text{right}\}$ , where  $\bar{y}_{in} = \arg \min_y l_{in}(y)$  is a “pseudolabel” indicating a child which gives the lowest value of the regression loss  $L$  for instance  $\mathbf{x}_n$  under the current tree.

For hyperplane nodes, this is NP-hard, but can be approximated by using a convex surrogate loss (we use the logistic loss). Hence, if  $\phi_i$  is an  $\ell_1$  norm, this requires solving an  $\ell_1$ -regularized logistic regression.

# Pseudocode for training a single TAO tree

TAO repeatedly alternates optimizing over sets of nodes by training a (binary) classifier in the decision nodes and a regressor in the leaves, while monotonically decr. the obj. function  $E(\Theta)$ .

```
input training set; initial tree  $\mathbf{T}(\cdot; \Theta)$  of depth  $\Delta$   
 $\mathcal{N}_0, \dots, \mathcal{N}_\Delta \leftarrow$  nodes at depth  $0, \dots, \Delta$ , respectively  
 $\mathcal{R}_1 \leftarrow \{1, \dots, N\}$   
repeat  
  for  $d = 0$  to  $\Delta$   
    parfor  $i \in \mathcal{N}_d$   
      if  $i$  is a leaf then  
         $\theta_i \leftarrow$  train regressor  $\mathbf{g}_i$  on reduced set  $\mathcal{R}_i$   
      else  
         $\theta_i \leftarrow$  train decision function  $f_i$  on  $\mathcal{R}_i$   
        compute the reduced sets of each child of  $i$   
until stop  
prune dead subtrees of  $\mathbf{T}$   
return  $\mathbf{T}$ 
```

# Ensemble of TAO trees

- Out of many ways to ensemble individual learners (bagging, boosting, etc.), we choose a simple one: we train each TAO tree independently (and in parallel) on a random subset of  $M$  samples of the available training data ( $N$  instances). If  $M = N$  this is bagging.

The forest prediction is the average of its trees' predictions.

- Although our TAO regression algorithm works more generally, our experiments use:
  - oblique trees ( $\mathbf{f}_i$  is linear), which are more powerful than axis-aligned trees
  - constant- and linear-predictor leaves ( $\mathbf{g}_i$ ).
- We initialize each TAO tree ( $\mathbf{T}$ ) from a complete tree of depth  $\Delta$  and random node parameters.
- We train each tree with an  $\ell_1$  regularizer to achieve a more compact model: it encourages the weight vectors of individual nodes to be sparse and leads to more dead subtrees which can be pruned.

# Experiments: standard benchmarks and algorithms

**TAO-c**: our oblique trees with constant leaves,

**TAO-l**: our oblique trees with linear leaves.

See the paper for extended results, additional datasets, etc.

	Forest	$E_{\text{test}}$	$T$	$\Delta$		Forest	$E_{\text{test}}$	$T$	$\Delta$
cpuct ( $N=8k, D=21$ )	CART	$3.63\pm 0.32$	1	25	CT slice ( $N=54k, D=384$ )	CART	$2.71\pm 0.06$	1	51
	<b>TAO-c</b>	$2.71\pm 0.04$	1	6		<b>TAO-c</b>	$1.54\pm 0.05$	1	7
	RF	$2.62\pm 0.04$	100	36		AdaBoost	$1.48\pm 0.03$	100	10
	ARF	$2.62\pm 0.01$	50	15		XGBoost	$1.45\pm 0.00$	100	10
	AdaBoost	$2.61\pm 0.16$	100	10		AdaBoost	$1.31\pm 0.01$	1k	10
	RF	$2.60\pm 0.01$	1k	37		XGBoost	$1.18\pm 0.00$	1k	10
	XGBoost	$2.60\pm 0.00$	100	10		<b>TAO-l</b>	<b><math>1.16\pm 0.02</math></b>	<b>1</b>	5
	ET	$2.58\pm 0.03$	100	45		ET	$1.06\pm 0.01$	100	82
	<b>TAO-l</b>	<b><math>2.58\pm 0.02</math></b>	<b>1</b>	5		RF	$1.03\pm 0.01$	100	71
	XGBoost	$2.57\pm 0.00$	1k	10		cRF	1.00	1k	-
	AdaBoost	$2.56\pm 0.11$	1k	10		RF	$0.97\pm 0.01$	1k	78
ET	$2.49\pm 0.03$	1k	50	<b>TAO-c</b>	$0.89\pm 0.02$	30	7		
<b>TAO-c</b>	$2.39\pm 0.05$	30	7	<b>TAO-l</b>	<b><math>0.58\pm 0.03</math></b>	<b>30</b>	6		
<b>TAO-l</b>	<b><math>2.35\pm 0.01</math></b>	<b>30</b>	5						

The **TAO** regression forests are smaller (fewer and shallower trees) yet consistently more accurate, particularly if using linear predictors at the leaves.

# Experiments: MNIST digit rotation

Task: given an MNIST image, predict a class-dependent rotation of it ( $N=60k, D=784, K=784$ ).

Forest	$E_{\text{test}} \times 10^{-2}$	#pars.	FLOPS	$T$	$\Delta$
AdaBoost	>24 hours runtime			39k	25
CART	$23.08 \pm 0.12$	120k	28	1	28
RF	$14.38 \pm 0.23$	7.6M	(2 830)	100	39
RF	$14.08 \pm 0.25$	68M	(28k)	1k	40
ET	$13.83 \pm 0.12$	12M	(3 091)	100	35
TAO-c	$13.76 \pm 0.09$	9M	42k	30	29
ET	$13.72 \pm 0.13$	109M	(3 360)	1k	38
XGBoost	$10.35 \pm 0.00$	180M	(613k)	39k	25
TAO-l	$9.63 \pm 0.17$	288k	4 491	1	7
TAO-l	$6.59 \pm 0.11$	7.7M	126k	30	7

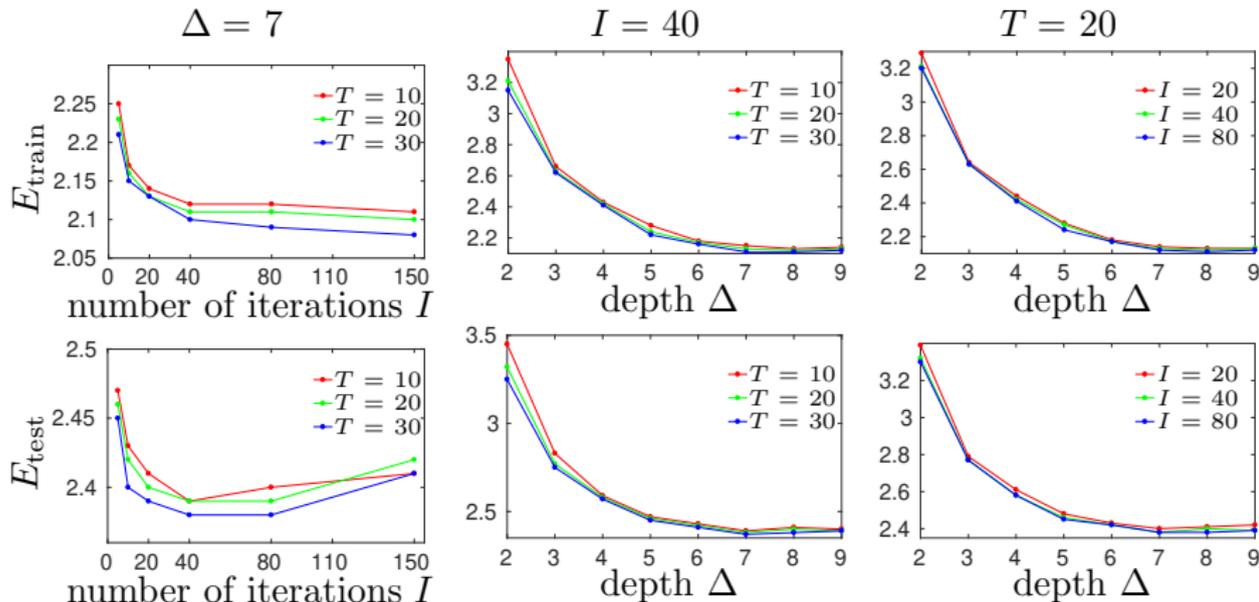
The improvement of TAO regression forests over other methods is clear, confirming our earlier results even more drastically.

# Experiments: MNIST digit rotation (cont.)

input										
ground-truth output										
RF $T = 1k$										
ET $T = 1k$										
XGBoost $T = 39k$										
TAO-1 $T = 1$										
TAO-1 $T = 30$										

# Experiments: forest hyperparameters

Hyperparameters' exploration (cpuact dataset): depth  $\Delta$ , # of TAO iterations  $I$  and # of trees  $T$ . Each column fixes one factor and varies the other two. In the paper, we also explore the effect of various diversity mechanisms.



Note how the forest can eventually overfit, which suggests that TAO is optimizing each tree well.

- Our hypothesis was that using more complex trees and a better tree optimization would produce more accurate forests. This was indeed the case, across a range of datasets we tried:
  - Our TAO regression forests outperform all competing algorithms we tested in terms of accuracy.
  - The TAO forests are smaller in terms of model size: number of trees, total number of parameters, depth.
- Their design in terms of hyperparameter tuning remains as simple as with random forests or boosting: we choose the tree depth and number of trees as large as computationally possible, but without overfitting.
- This makes our TAO forests a model of immediate, widespread practical applicability and impact, and suggests it could become the state-of-the-art in tree ensembles.

In separate papers, we have also found that TAO trees improve significantly in classification (rather than regression) and with boosting (rather than bagging).