

TOWARDS BETTER DECISION FORESTS: FOREST ALTERNATING OPTIMIZATION Miguel Á. Carreira-Perpiñán and Magzhan Gabidolla and Arman Zharmagambetov, EECS, UC Merced

1 INTRODUCTION

Decision forests (ensembles of decision trees) are widely recognized as amo the most accurate ML models for many tasks. However, neither the individ trees nor the forest are constructed to optimize a specific loss function.

In a series of papers, we have given algorithms that optimize very general typ of losses (in the sense of monotonically decreasing the loss over iterations) or a single tree (axis-aligned or oblique), in combination with popular ensembl mechanisms (bagging, AdaBoost, gradient boosting), and here over all the for parameters jointly, in all cases consistently improving over the state-of-the-(such as XGBoost or LightGBM).

Work partially supported by NSF award IIS–2007147.

2 SINGLE TREE: TREE ALTERNATING OPTIMIZATION (TAC

A scalable algorithm that can take a tree of arbitrary but parametric form monotonically decrease an objective function of the form loss + regularization

$$\min_{\boldsymbol{\tau}} \sum_{n} L(\mathbf{y}_n, \boldsymbol{\tau}(\mathbf{x}_n; \{\mathbf{w}_i\})) + \lambda \sum_{i \in \text{ nodes of } \boldsymbol{\tau}} \phi(\mathbf{w}_i)$$

We focus on oblique trees (which are far more powerful than axis-aligned ones

- Decision nodes: (sparse) hyperplane
- Leaf nodes: constant label or value

No gradient descent (the tree defines a piecewise constant function) but alterr ing optimization over the nodes. Based on two theorems:

- Separability condition: the objective function separates over nodes which a not descendant of each other.
- Reduced problem over a node: optimizing over a node's parameters takes special form that can be solved exactly or approximately:
- decision node: weighted 0/1 loss binary classification
- leaf node: majority vote or average

TAO operates on an initial tree structure. The final structure is usually a subse this because pruning occurs automatically via a ℓ_1 penalty on the decision nod weights: $\phi(\mathbf{w}_i) = \|\mathbf{w}_i\|_1$. This also sparsifies the decision hyperplanes.

The table below shows how a single TAO tree improves upon the traditional CA on test accuracy for several classification benchmarks:

| Algorithm | MNIST | Letter | SensIT | Pendigits | Spambase |
|-----------|-------|--------|--------|-----------|----------|
| CART | 88.05 | 86.07 | 81.00 | 91.62 | 89.62 |
| TAO | 94.74 | 90.41 | 85.44 | 96.80 | 93.31 |

Examples of tree-based models trainable with TAO:

- [1]: sparse oblique tree
- [4]: tree of neural nets
- [7]: softmax tree
- [8]: nonlinear embeddings with trees
- [12]: semi-supervised learning with trees
- [11]: clustering with trees

| | B F | FOREST: TA | O + BAG | GING/B | OOSTING | | | |
|-------------|--|------------------------------|--|----------------|---|------------------------------------|-----------------|------|
| nong | | | | | nsembling mech | anism resi | ults in bet | tter |
| idual | | s (higher accur | | - | | | | |
| | | gging [2], [3] | | | | | | |
| ypes | | | [10] | | | | | |
| over | AdaBoost [5], [6], [10] | | | | | | | |
| bling | gradient boosting [9] | | | | | | | |
| orest | We use sparse oblique trees, which are a much stronger learner than the tradi- | | | | | | | |
| e-art | | axis-aligned tre | | | | _ | | |
| c-ari | • mc | pre powerful mo | odel: obliqu | ie rather th | nan axis-aligned | l tree | | |
| | bet | tter optimizatio | n: TAO rath | er than C | ART/C5.0/etc. | | | |
| | This re | esults in more | accurate fo | rests; the | ensemble is div | erse enou | յh. | |
| | The ta | able below show | ws how sim | iple TAO ti | rees in bagging | outperform | other esta | ab- |
| \O) | lished | tree ensemble | es for regre | ssion prot | olems. | | | |
| and | | cpuact (| N=8k,D=21) | - | CT slice (N | V=54k, <i>D</i> =384 | .) | |
| n: | - | Forest | E _{test} | $T \Delta$ | Forest | E _{test} | $T \Delta$ | |
| | | CART | 3.63±0.32 | 1 25 | | 2.71±0.06 | 1 51 | |
| (1) | | Bagging-TAO Random Forest | $2.71{\pm}0.04$ $2.62{\pm}0.04$ | 1 6 100 36 | Bagging-TAO AdaBoost | $1.54{\pm}0.05$ $1.48{\pm}0.03$ | 100 10 | |
| ς. | | AdaBoost | 2.61 ± 0.04 | 100 10 | XGBoost | 1.45 ± 0.00 | 100 10 | |
| es): | | Random Forest | $2.60{\pm}0.01$ | 1k 37 | AdaBoost | $1.31{\pm}0.01$ | 1k 10 | |
| | | XGBoost | 2.60 ± 0.00 | 100 10 | XGBoost | 1.18 ± 0.00 | 1k 10 | |
| | | XGBoost AdaBoost | $2.57{\pm}0.00$ $2.56{\pm}0.11$ | 1k 10 1k 10 | Random Forest Random Forest | $1.03{\pm}0.01 \\ 0.97{\pm}0.01$ | 100 71 1k 78 | |
| rnat- | | Bagging-TAO | 2.39 ± 0.05 | 30 7 | Bagging-TAO | 0.89 ± 0.01 | 30 7 | |
| | - | | | | | | | |
| are | | 28 | | | 28 · hr | | - | |
| | | 26 24 | B-sklearn |] | 26 | | | |
| es a | | 22 | XG | Boost | 24 Kann | GB-skle | arn | |
| | | 20 - | and the second sec | - | $\begin{array}{c} 22 \\ 20 \end{array}$ | it man | ~~~~ | |
| | | | | ThtGBM | 20 - 18 - | LightG | BM | |
| | | 16 - | GB-TAO | | 16 GB-TA | | | |
| set of | | 0 2 | 4 6 | 8 | | 150 200 250 | 0 300 | |
| odes' | | | r of parame | | | ing steps 250 | | |
| | | Fiq | jure: Compar | rison of metl | hods for news20 da | itaset. | | |
| ART | Howe | ver. while each | individual | tree is we | Il optimized, the | forest is n | ot. Trees a | are |
| | | d independently | | | | | | • |
| | | | (| | J (| | | |

References:

- [1] Alternating optimization of decision trees, with application to learning sparse oblique trees, NEURI
- [2] Smaller, more accurate regression forests using tree alternating optimization. *ICML 2020*.
- [3] Ensembles of bagged TAO trees consistently improve over random forests, AdaBoost and gradient FODS 2020.
- [4] Learning a tree of neural nets. *ICASSP 21*.
- [5] Improved boosted regression forests through non-greedy tree optimization. *IJCNN 21*.
- [6] Improved multiclass AdaBoost for image classification: The role of tree optimization. *ICIP 2021*.
- [7] Softmax tree: An accurate, fast classifier when the number of classes is large. *EMNLP 2021*.
- [8] Learning interpretable, tree-based projection mappings for nonlinear embeddings. *AISTATS 2022*.
- [9] Pushing the envelope of gradient boosting forests via globally-optimized oblique trees. CVPR 22.
- [10] Improved multiclass AdaBoost using sparse oblique decision trees. *IJCNN 22*.
- [11] Optimal interpretable clustering using oblique decision trees. *KDD 2022*.
- [12] Semi-supervised learning with decision trees: Graph Laplacian tree alternating optimization. NEUF

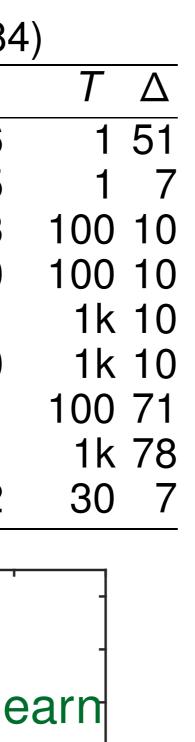
4 FOREST ALTERNATING OPTIMIZATION (FAO)

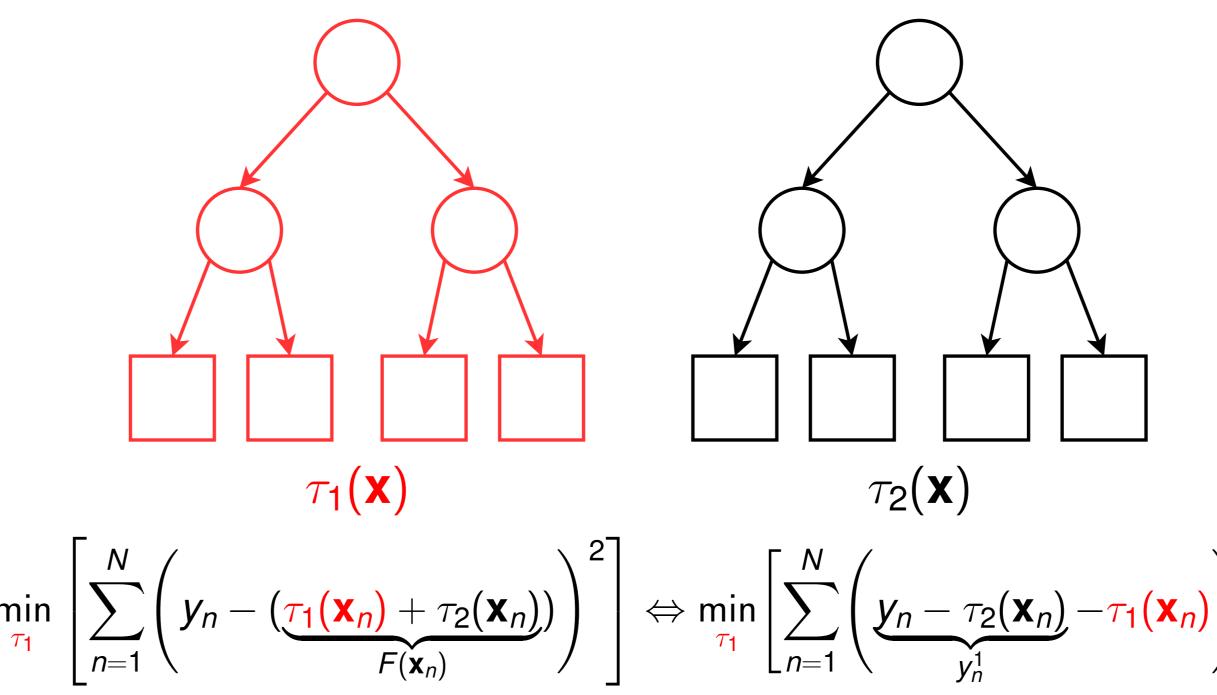
In this work, we take this one step further and optimize globally over all the parameters (decision & leaf nodes) of a forest having a fixed number of trees of given structure, monotonically decreasing an objective function of the form loss + regularization:

 $\min_{\boldsymbol{\tau}_1,\ldots,\boldsymbol{\tau}_T}\sum L(\mathbf{y}_n,$

where $\mathbf{F}(\mathbf{x}) = \sum_{t=1}^{T} \boldsymbol{\tau}_t(\mathbf{x})$ is a forest of T trees. Alternating optimization over trees:

• If we fix all trees but one, the resulting problem over that tree can be optimized by TAO.

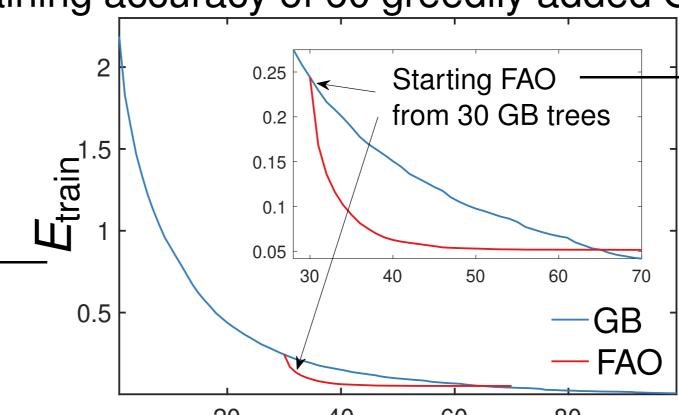




FAO is very good at optimizing the model. The plot shows how 30 trees trained with FAO can exceed the training accuracy of 60 greedily added GB trees.

However, the resulting forest can easily overfit: the plot shows the increasing test error of FAC trained on 10 trees of various initialization.

Щ²⁺



Tree additions/FAO iterations

Some experiment results:

| RIPS 2018. | MNIST (<i>N</i> = 60 <i>k</i> , <i>D</i> = 784, <i>K</i> = 10) | | | | | |
|--------------|---|-------------------|--------|------|----|----------|
| nt haasting | Forest | Etest (%) | #pars. | Т | Δ | Forest |
| nt boosting. | SPORF | 2.89±0.04 | (143M) | 1k | 50 | SPORF |
| | XGBoost | $2.20{\pm}0.00$ | 107k | 1k | 6 | SPORF |
| | LightGBM | $2.02{\pm}0.00$ | 121k | 1k | 10 | XGBoost |
| | XGBoost | $1.91 {\pm} 0.00$ | 505k | 10k | 6 | XGBoost |
| | GB-TAO | $1.65{\pm}0.02$ | 3M | 500 | 7 | LightGBM |
| 2. | LightGBM | $1.62{\pm}0.00$ | 642k | 10k | 21 | LightGBM |
| | GB-TAO | $1.55{\pm}0.02$ | 7.2M | 1.4k | 7 | XGBoost |
| | avg-FAO | $1.48{\pm}0.06$ | 658k | 60 | 6 | LightGBM |
| | avg-FAO | $1.39{\pm}0.04$ | 968k | 90 | 6 | avg-FAO |
| URIPS 2022. | avg-FAO | $1.33{\pm}0.04$ | 4.9M | 300 | 8 | avg-FAO |
| | | | | | | |



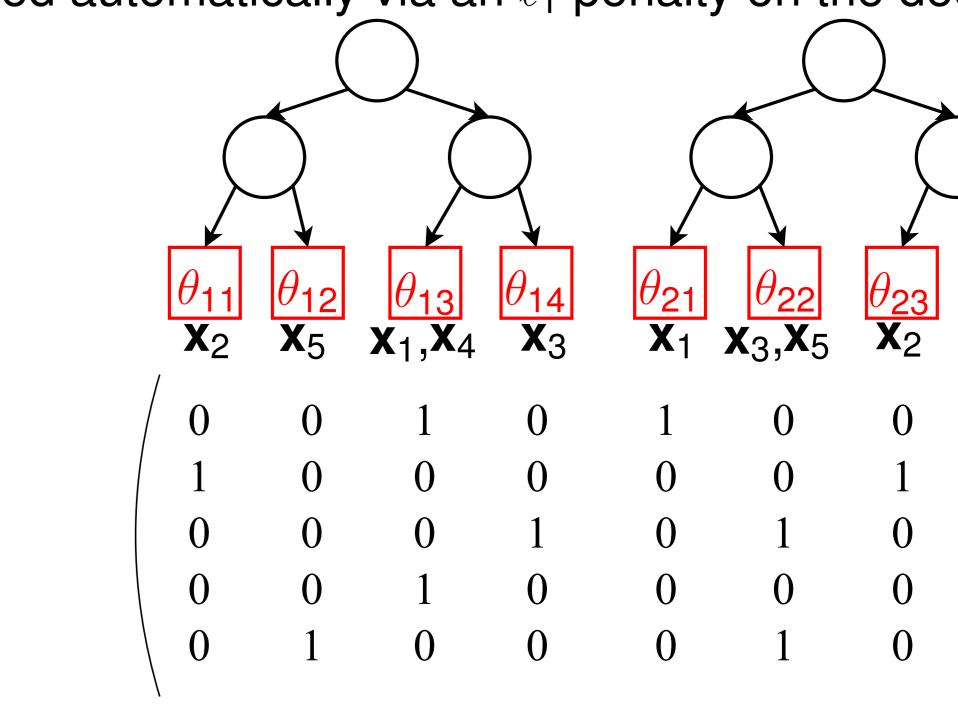


$$\mathbf{F}(\mathbf{x}_n)) + \lambda \sum_{t=1}^{\prime} \sum_{i \in \text{ nodes of } \boldsymbol{\tau}_t} \phi(\mathbf{w}_{ti})$$

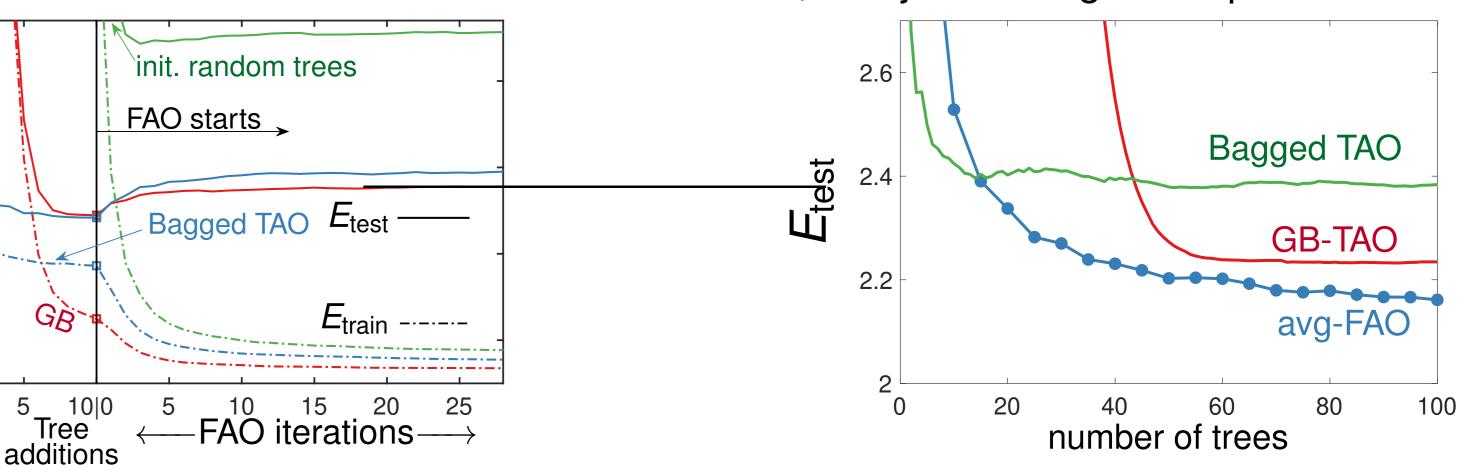
(2)

*Y*5

• Also, if we fix all the decision nodes of all the trees, the resulting problem over all the trees' leaves can be optimized exactly. As with a single tree, each individual tree's structure is still pruned automatically via an ℓ_1 penalty on the decision nodes' weights.



To obtain better generalization we train multiple small FAO forests on different random initializations, and just average their predictions.



Y (N = 4.5M, D = 18, K = 2) E_{test} (%) #pars. (271M) 100 19.91 102 19.73 (2.7B) 1k 109 151k 300 19.63 19.62 196k 100 19.62 153k 100 19.60 230k 300 19.59 2.0M 1k 10 1.5M 19.57 1k 233k 19.5 50 19.50 459k 100

casp (N = 45k, D = 9) $E_{\text{test_rmse}}$ #pars. $T \Delta$ Forest 119k 100 10 XGBoost $3.66 {\pm} 0.00$ 793k 1k 10 $3.58{\pm}0.00$ XGBoost LightGBM 3.54±0.00 153k 100 114 GB-TAO 256k 50 12 $3.49{\pm}0.01$ LightGBM 3.48±0.00 766k 1k 109 avg-FAO 359k 50 12 3.45 ± 0.02 GB-TAO 481k 100 12 $3.43{\pm}0.00$ 711k 100 12 $3.40{\pm}0.01$ avg-FAO GB-TAO $3.39{\pm}0.01$ 887k 200 12 1.4M 200 12 avg-FAO $3.37{\pm}0.01$