

Convex Polytope Trees: Appendix

A Proofs

In this section, we show why the splitting function at each internal node results in a convex set confined within a convex polytope. We start by proving a lemma which is needed to prove the main theorem.

Lemma 1. For any $\{r_i, \beta_i\}_{i=1}^K$, such that $r_i \in \mathbb{R}_+$ and $\beta_i \in \mathbb{R}^d$, the function:

$$g(\mathbf{x}) := \sum_{i=1}^K r_i \ln(1 + e^{\beta_i' \mathbf{x}}) \quad (18)$$

is convex over its domain \mathbb{R}^d .

Proof. Since the sum of convex functions is also convex, it suffice to show each term of g is a convex function. We demonstrate this by using the following theorem: ‘‘A function is convex iff its second derivative is a positive semi-definite matrix over the domain.’’ One can omit r_i ’s in the following calculations because a positive scalar does not change the convexity.

The first derivative of each term is:

$$\frac{\partial \ln(1 + e^{\beta_i' \mathbf{x}})}{\partial \mathbf{x}} = \frac{e^{\beta_i' \mathbf{x}}}{e^{\beta_i' \mathbf{x}} + 1} \beta_i' \quad (19)$$

and by taking the derivative of the above vector, we will have:

$$\frac{\partial^2 \ln(1 + e^{\beta_i' \mathbf{x}})}{\partial \mathbf{x}^2} = \frac{e^{\beta_i' \mathbf{x}}}{(e^{\beta_i' \mathbf{x}} + 1)^2} \beta_i \beta_i' \quad (20)$$

where $\beta_i \beta_i'$ is a matrix in $\mathbb{R}^d \times \mathbb{R}^d$. Since the scalar $\frac{e^{\beta_i' \mathbf{x}}}{(e^{\beta_i' \mathbf{x}} + 1)^2}$ is positive for any \mathbf{x} , we just need to show the matrix $\beta_i \beta_i'$ is positive semi-definite. To that end, we prove for any $\mathbf{v} \in \mathbb{R}^d$:

$$\mathbf{v}' (\beta_i \beta_i') \mathbf{v} \geq 0.$$

And, that can be shown by:

$$\mathbf{v}' (\beta_i \beta_i') \mathbf{v} = (\mathbf{v}' \beta_i) \cdot (\beta_i' \mathbf{v}) = (\beta_i' \mathbf{v})' \cdot (\beta_i' \mathbf{v}) = \|\beta_i' \mathbf{v}\|^2 \geq 0.$$

Therefore the proof of the lemma is complete. □

Theorem 1. For any $\{r_i, \beta_i\}_{i=1}^K$, such that $r_i \in \mathbb{R}_+$ and $\beta_i \in \mathbb{R}^d$, let:

$$A_{left} = \{\mathbf{x} \mid \mathbf{x} \in \mathbb{R}^d, f(\mathbf{x}) \leq q_{thr}\}, \quad \text{where: } f(\mathbf{x}) = 1 - e^{-\sum_{i=1}^K r_i \ln(1 + e^{\beta_i' \mathbf{x}})} \quad (4)$$

then A_{left} is a convex set, confined by a convex polytope.

Proof. We start by showing A_{left} is a convex set. Note that, due to the duality

$$\mathbf{x} \in A_{left} \iff f(\mathbf{x}) \leq q_{thr} \quad (21)$$

By the definition of a convex set, we just need to prove the following:

$$\forall t \in [0, 1], \forall \mathbf{x}_1, \mathbf{x}_2 \in \mathbb{R}^d \quad \text{if } f(\mathbf{x}_1), f(\mathbf{x}_2) \leq q_{thr} \implies f(t\mathbf{x}_1 + (1-t)\mathbf{x}_2) \leq q_{thr}. \quad (22)$$

Let $g(\cdot)$ and q_{thr}^* be:

$$g(\mathbf{x}) := -\ln(1 - f(\mathbf{x})) = \sum_{i=1}^K r_i \ln(1 + e^{\beta_i' \mathbf{x}}), \quad q_{thr}^* := -\ln(1 - q_{thr}). \quad (23)$$

Since $-\ln(1-a)$ is montically increasing with respect to a , replacing f by $g(\cdot)$ and q_{thr} by q_{thr}^* in (22), results in a mathematically equivalent expression. Now, we can prove the new statement using Jensen’s inequality. To be more specific, based on Lemma 1 (g is convex) and Jensen’s inequality, we have:

$$\forall t \in [0, 1], \forall \mathbf{x}_1, \mathbf{x}_2 \in \mathbb{R}^d \quad g(t\mathbf{x}_1 + (1-t)\mathbf{x}_2) \leq tg(\mathbf{x}_1) + (1-t)g(\mathbf{x}_2). \quad (24)$$

So if $g(\mathbf{x}_1), g(\mathbf{x}_2) \leq q_{thr}^*$:

$$g(t\mathbf{x}_1 + (1-t)\mathbf{x}_2) \leq tg(\mathbf{x}_1) + (1-t)g(\mathbf{x}_2) \leq tq_{thr}^* + (1-t)q_{thr}^* = q_{thr}^*$$

proving A_{left} is convex.

We are just remained with showing A_{left} is confined within a convex polytope. This can be shown by:

$$\begin{aligned} f(\mathbf{x}) \leq q_{thr} &\iff g(\mathbf{x}) \leq q_{thr}^* \implies \forall i \in [1 : K], \quad r_i \ln(1 + e^{\beta_i' \mathbf{x}}) \leq q_{thr}^* \\ &\iff \forall i \in [1 : K], \quad \beta_i' \mathbf{x} \leq \ln(e^{\frac{q_{thr}^*}{r_i}} - 1) \end{aligned}$$

which completes the proof. \square

B Additional details on experimental settings

As mention in the paper, we train CPT in a probabilistic manner and switch to a deterministic tree at test time. To make the transition smoother, we conduct annealing during training. To be more specific, we transform the probability function $f(\mathbf{x})$ at each node to $f_{\lambda_t}(\mathbf{x})$, where:

$$f(\mathbf{x}) = 1 - e^{-\sum_{i=1}^K r_i \ln(1 + e^{\beta_i' \mathbf{x}})} \quad \text{and} \quad f_{\lambda_t}(\mathbf{x}) := \frac{1}{1 + \left(\frac{1-f(\mathbf{x})}{1-p_0}\right)^{\lambda_t}} \quad (25)$$

Larger λ_t results in a sharper change of probability from 0 to 1, and p_0 controls where that change happens. During training, we gradually increase λ_t to make the gap between probabilistic and deterministic tree progressively smaller. We also learn p_0 like other parameters of the model using SGD. Notice, the change of f to f_{λ_t} keeps the mathematical and geometrical interpretation of CPT intact. That is because any thresholding of f_{λ_t} has an equivalent counterpart for f since f and f_{λ_t} are strictly monotonic function of each other.

Note, the main goal of our comparison in Table 2 is reporting the highest possible value a method can achieve regardless of the size. This has significant importance in the decision tree literature, as they are often undermined since their accuracy is lower than the NN counterparts, and it is important to push the accuracy limit of these classes of methods because they provide interpretability which is missing with NNs.

For the comparison, we tried our best to state the best performance reported for each method and if there were no results on a specific dataset and code was available, we trained their model up to the depth of 16. For the FTEM method, the reported numbers are based on the author’s best hyper-parameter tuning. They varied the maximum depth from 2-18 with a step size of 2. These numbers were only available for MNIST, Connect4, and SensIT. For other classification results, we trained their model and reported the results. However, their method was not applicable to regression, so we did not report the results.

For TAO, the code is not available. However, the authors of TAO have written a follow-up paper [50] comparing old and new decision trees. We use the results in this paper to report the performance of TAO on MNIST, Connect4, SensIT, and Letter. They did not provide all the details of how they varied the maximum length, but for some methods, they varied the depth up to 30. And for the results on HIV and Bace, we relied on the reported numbers in LCN (Lee and Jaakola, 2020), where they vary the depth of the tree in the interval [2, 12].

For the LCN method on HIV, Bace, and PDBbind, we used their reported results in which they varied the depth from 2, 3, ..., 12. It is worth noting that LCN does not learn the tree structure, and it always learns a full binary tree.

We did not compare our method with soft trees (soft at test time) that are often not considered interpretable. For this reason, they have not been previously compared to the deterministic trees in

the literature. However, it is worth mentioning that CPT results for the soft version considerably improve if we use it as a soft tree at the test time. For instance on the Bace and HIV datasets, we can achieve 1-2 percentage increase in AUC if we use the soft version. The only closely related method is FTEM. During the training period, they are also a probabilistic tree and use the thresholded soft tree at the test time, which we do compare against in our experiments.

We performed the classification and regression experiments on a laptop with a 2.5 GHz 6-Core Intel Core i7 CPU and 16 GB of RAM.