

Binary Interpolation Tree (B.I.T)

Binary Interpolation Tree is a machine learning model i have been working on mostly out of curiosity. Please let me know if my terminology is off, i made any other errors, or something very similar already exists.

The main idea was to conceptualize a model that provides:

- fast runtime on cpu's in relation to model complexity ($O(n)$, $\Omega(1)$, input dependent)
- the ability to further simplify computations based on a maximum inaccuracy
- simple and fast training
- very low initial model complexity that increases during training
- multithreading capabilities for maximum cpu usage
- increased transparency compared to other ML models
- a parameter regarding determinism of the algorithm

1. The Model

1.1 The basic model

As the name suggests, this model is basically a binary tree. This binary tree consists of nodes and end-nodes. Nodes always have 2 distinct children except for end-nodes that have none. Every node has a reference to its parent, except for the root-node, which has no parent. Every node has a position vector and every end-node additionally contains an output vector. For the purpose of describing this system mathematically, every node has an index i and the following variables associated with it:

The index of the first child: $c_1(i)$.

The index of the second child: $c_2(i)$.

The index of the parent: $parent(i)$.

The position vector: $\vec{p}(i) \in V_i$.

The result of the calculation for nodes, or the output vector for end-nodes: $\vec{o}(i) \in V_o$.

For the root-node the index $i = 0$.

An index of -1 means there is no reference. For example $parent(0) := -1$.

To calculate the output $\vec{o}(0)$ of this model for a given input vector $\vec{x} \in V_i$, this binary tree is traversed recursively in my cpu implementation, from the root-node, to the end-nodes. The following calculations are performed for every node:

$$I(\vec{p}_1, \vec{p}_2, \vec{x}) := \frac{(\vec{p}_2 - \vec{p}_1) \circ (\vec{x} - \vec{p}_1)}{\|\vec{p}_2 - \vec{p}_1\|^2}$$

$$clamp(x) := \min(1, \max(0, x))$$

$$a(c_1(i)) := \text{clamp}(I(\vec{p}(c_1(i)), \vec{p}(c_2(i)), \vec{x})) \quad a(c_2(i)) := 1 - a(c_1(i))$$

$$\vec{o}(i) := a(c_1(i)) \cdot \vec{o}(c_1(i)) + a(c_2(i)) \cdot \vec{o}(c_2(i))$$

if $a(c_1(i)) = 0$ or $a(c_2(i)) = 0$, $\vec{o}(c_1(i))$ or $\vec{o}(c_2(i))$ doesn't have to be calculated, which can cause a great performance increase for sequential execution. the following statements are true and fully define $I(\vec{p}_1, \vec{p}_2, \vec{x})$ in an intuitive way:

$$I(\vec{p}_1, \vec{p}_2, a * \vec{p}_1 + (1 - a) * \vec{p}_2) = 1 - a \Rightarrow I(\vec{p}_1, \vec{p}_2, \vec{p}_1) = 0 \wedge I(\vec{p}_1, \vec{p}_2, \vec{p}_2) = 1$$

$$I(\vec{p}_1, \vec{p}_2, \vec{x}_1) = I(\vec{p}_1, \vec{p}_2, \vec{x}_2) \Leftrightarrow (\vec{p}_2 - \vec{p}_1) \circ (\vec{x}_1 - \vec{x}_2) = 0.$$

1.2 Introducing indeterministic behavior

Swapping the definition of $a(c_1(i))$ from $a(c_1(i)) := \text{clamp}(I(\vec{p}(c_1(i)), \vec{p}(c_2(i)), \vec{x}))$ to

$$a(c_1(i)) := \text{clamp}(I(\vec{p}(c_1(i)), \vec{p}(c_2(i)), \vec{x}) * d + (1 - d) * r)$$

where $0 \leq d \leq 1 \wedge 0 \leq r \leq 1$, will allow indeterministic behavior, if a random number is assigned to r for each computation of $a(c_1(i))$. The parameter d is called determinism in my implementation, because $d = 1$ means that the model is 100% deterministic while $d = 0$ means the interpolation factor $a(c_1(i))$ is completely random between 0 and 1.

1.3 Simplifying computations by defining a maximum inaccuracy

For this purpose, the following variables are defined:

The number of dimensions of the input space: $m := \dim(V_i)$.

The number of dimensions of the output space: $n := \dim(V_o)$.

The Inaccuracy-scale vector $\vec{s} \in V_o$ allows different maximum inaccuracies for all different output dimensions and is meant to be immutable for every model. The maximum inaccuracy t is meant to be passed for every Forward pass, to be able to calculate the output for the same model with varying speed and accuracy. This step profits from normalizing the Data beforehand. The following variables are defined for each node:

The complexity: $\text{complexity}(i) := 1 + \text{complexity}(c_1(i)) + \text{complexity}(c_2(i))$ for all nodes and $\text{complexity}(i) := 1$ for all end-nodes specifically.

$\text{complexity}(i)$ is roughly proportional to the runtime of the forward pass off the node with the index i .

$$\text{The contribution to the final result: } c(i) := \begin{cases} c(\text{parent}(i)) * a(i) & i \neq 0 \\ 1 & i = 0 \end{cases}$$

$$\text{The maximum contribution value } v(i) := \max_{1 \leq j \leq n} \left(\frac{\vec{o}(i)_j}{\vec{s}_j} \right) \text{ for all end-nodes and}$$

$v(i) := \max(v(c_1(i)), v(c_2(i)))$ for every node that is no end-node. $v(i)$ can be cached for every node, because \vec{s} is meant to be immutable for every model.

$$\sum_{k \in M} v(k) \cdot c(k) \leq t \Rightarrow \forall 1 \leq j \leq n : \sum_{k \in M} \vec{o}(k)_j * c(k) \leq \vec{s}_j * t$$

This Expression means, that you can set $a(k) := 0$ for $k \in M$ and can guarantee that the inaccuracy of $\vec{o}(0)$, caused by this simplification, will be less than $\vec{s} * t$ for each output dimension, which is a way of trading accuracy for speed. In my c# implementation, i wrote the recursive forward pass as an iterator method, to collect all nodes where $v(i) \cdot c(i) \leq t$ (potential minors) and choose as much nodes as possible,

that have the the lowest $\frac{v(i) \cdot c(i)}{\text{complexity}(i)}$ ratio from those potential minors to be in M , so that

$\sum_{k \in M} v(k) \cdot c(k) \leq t$ is still satisfied. All nodes in M are called minors and get selected by

$\frac{v(i) \cdot c(i)}{\text{complexity}(i)}$ because of the following statements:

$v(i_1) \cdot c(i_1) < v(i_2) \cdot c(i_2) \wedge \text{complexity}(i_1) = \text{complexity}(i_2)$ means that it is beneficial to choose the node with the index i_1 as a minor, because this introduces less in inaccuracy, but roughly the same runtime.

$v(i_1) \cdot c(i_1) = v(i_2) \cdot c(i_2) \wedge \text{complexity}(i_1) < \text{complexity}(i_2)$ means that it is beneficial to choose the node with the index i_1 as a minor, because this introduces the same inaccuracy, but it tends run the forward pass faster.

$v(i_1) \cdot c(i_1) + v(i_2) \cdot c(i_2) = v(i_3) \cdot c(i_3) \wedge \text{complexity}(i_1) + \text{complexity}(i_2) = \text{complexity}(i_3)$ means that choosing the nodes with the index i_1 and i_2 as a minor, will have almost the same impact as choosing only the node with the index i_3 , because this introduces the same inaccuracy, and roughly the same runtime.

2. Supervised training of the model

2.1 Training the parameters

For each input vector $\vec{x}(j)$ where $1 \leq j \leq b$ and b is the size of the batch, the result of the forward pass $\vec{o}_j(0)$ is calculated. $\vec{y}(j)$ represents the target $\vec{o}_j(0)$ is meant to approximate for the given input vector $\vec{x}(j)$. For every node that is no end-node, the index j is necessary for $\vec{o}_j(i)$, because $\vec{o}_j(i)$ depends on $\vec{x}(j)$, but for every end-node, $\vec{o}(i)$ remains constant within one training iteration and is therefore redundant. The function $error(\vec{w}_1, \vec{w}_2)$ represents the error metric and $c_j(i)$ is the contribution value for each node with the index i and input vector $\vec{x}(j)$. The set $S(i)$ contains all indices j for each node with the index i , where $c_j(i) \neq 0$. After each training iteration, $\alpha \cdot \Delta \vec{p}(i)$ is added to $\vec{p}(i)$, where α is the learning rate of this training process. Intuitively speaking, the position vector $\vec{p}(i)$ is adjusted towards those $\vec{x}(j)$ with $j \in S(i)$, the node is most „responsible“ for and where the biggest error occurs. If $|S(i)| = 0$, all parameters that belong to the node with the index i remains unchanged.

$$\Delta \vec{p}(i) := \frac{1}{|S(i)|} \sum_{j \in S(i)} (\vec{x}(j) - \vec{p}(i)) * c_j(i) * \left(1 - \exp\left(-\frac{error(\vec{o}_j(0), \vec{y}(j))}{f(i)}\right) \right)$$

$$f(i) := \frac{1}{|S(i)|} \sum_{j \in S(i)} error(\vec{o}_j(0), \vec{y}(j))$$

When the position vector of all nodes was adjusted, the only parameters left unchanged are those within the output vector of each end-node. To adjust those, $\alpha \cdot \Delta \vec{o}(i)$ is added to $\vec{o}(i)$.

$$\Delta \vec{o}(i) := \frac{1}{|S(i)|} \sum_{j \in S(i)} (\vec{y}(j) - \vec{o}_j(0)) * c_j(i)$$

The preceding formulas also show that $\vec{o}_j(0)$ can be modified with functions like softmax for example, where the first derivative is always positive. This comes with the cost of possibly „overshooting“ the target by adding $\alpha \cdot \Delta \vec{o}(i)$ to $\vec{o}(i)$ and divergence, which causes the training to fail if alpha is too large.

2.2 Restructuring the binary tree

This model has the ability to begin its training with the lowest complexity possible that does not output a constant: a root-node, that has 2 end-nodes as its children. Over time, the complexity of the model can be increased with little computational effort and no negative impact on its accuracy. For this purpose, the standard deviation is used to determine which branches of this binary tree require the most complexity. Intuitively speaking, $(\vec{y}(j) - \vec{o}_j(0)) * c_j(i)$ ideally shouldn't deviate from $\alpha \cdot \Delta \vec{o}(i)$ for all $j \in S(i)$, because this means that this model can depict the Algorithm $f(\vec{x})$ it is meant to approximate perfectly. If it does deviate however, this branch is either not complex enough, meaning that $complexity(i)$ should be greater than it is, or $f(\vec{x})$ is nondeterministic and can not be depicted perfectly. If this deviation is small for any node, both children of this node are end-nodes and have very similar output vectors, this node can be replaced by an end-node that has the same position vector as this node and the mean of the output vectors of its children as the output vector, This way redundant complexity can be reduced, without $\vec{o}(0)$ changing to much for any input as a result.