

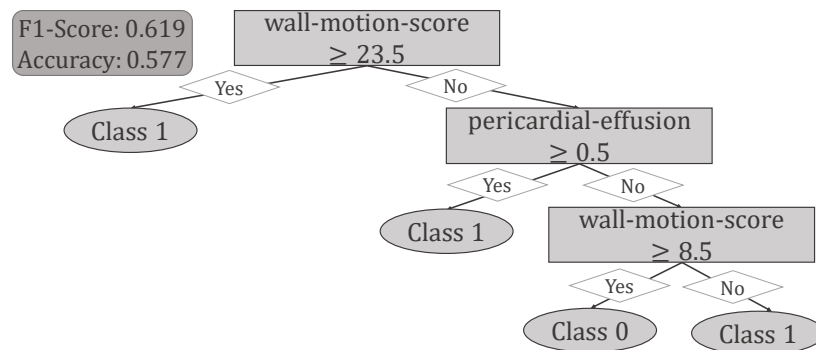
GradTree: Learning Axis-Aligned Decision Trees with Gradient Descent



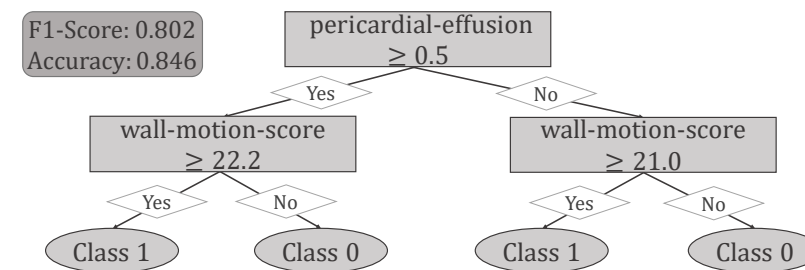
Motivation

Decision Trees (DT) are commonly used for many machine learning tasks due to their high interpretability

- Learning a DT is a difficult optimization problem (non-convex and non-differentiable)
- Prevailing approach is a greedy procedure, minimizing the objective locally at each internal node
 - Constrains search space and potentially leads to suboptimal trees



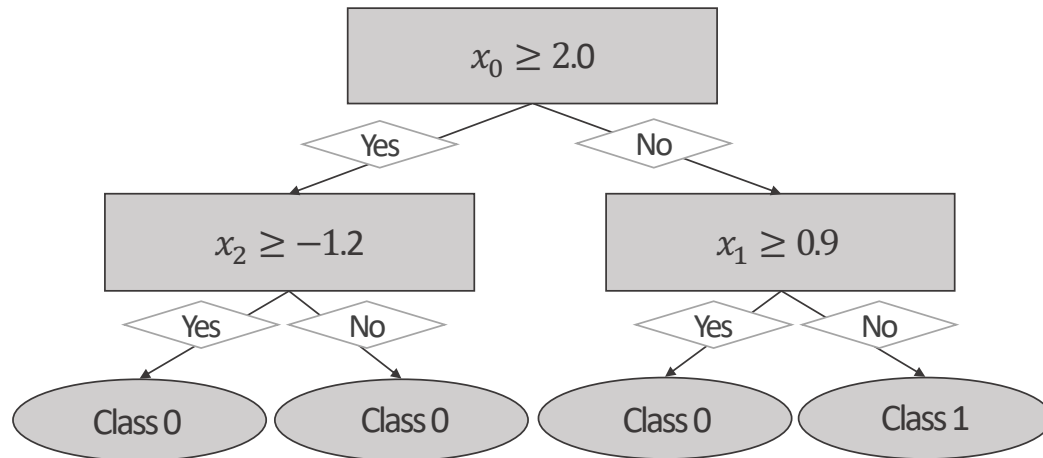
(a) Greedy DT (CART)



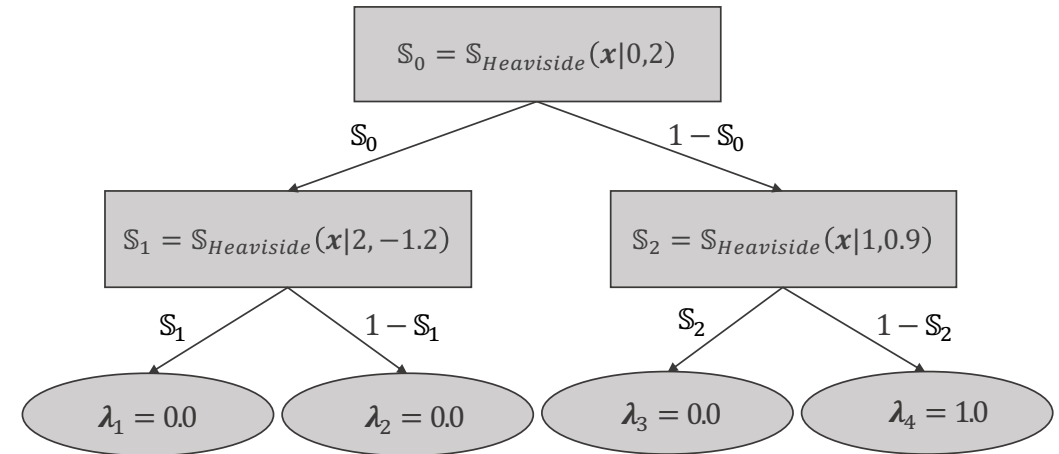
(b) Gradient-Based DT (GradTree)

Figure 1: **Greedy vs. Gradient-Based DT.** Two DTs trained on the Echocardiogram dataset. The CART DT (left) makes only locally optimal splits, while GradTree (right) jointly optimizes all parameters, leading to significantly better performance.

Arithmetic Decision Tree Formulation



1. If $(x_0 \geq 2.0) \wedge (x_2 \geq -1.2)$, then assign to Class 0.
2. If $(x_0 \geq 2.0) \wedge (x_2 \leq -1.2)$, then assign to Class 0.
3. If $(x_0 \leq 2.0) \wedge (x_1 \geq 0.9)$, then assign to Class 0.
4. If $(x_0 \leq 2.0) \wedge (x_1 \leq 0.9)$, then assign to Class 1.



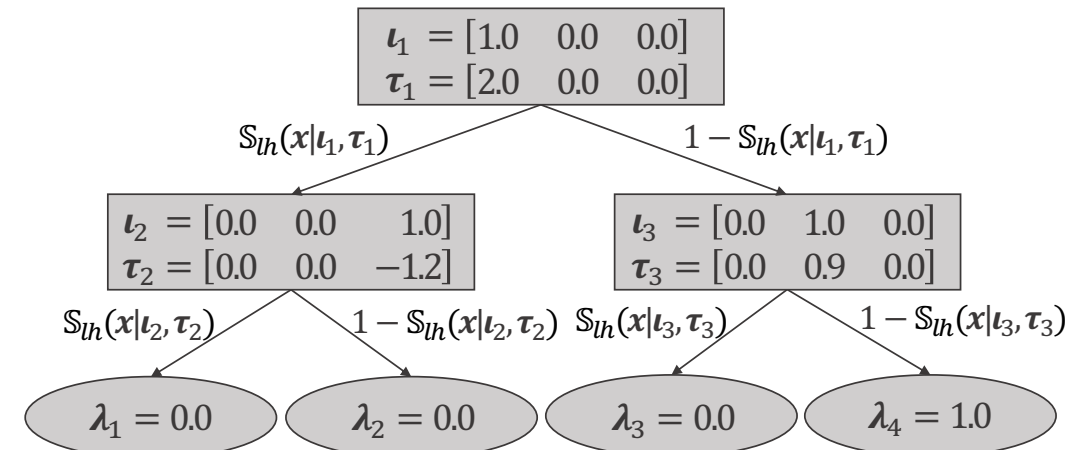
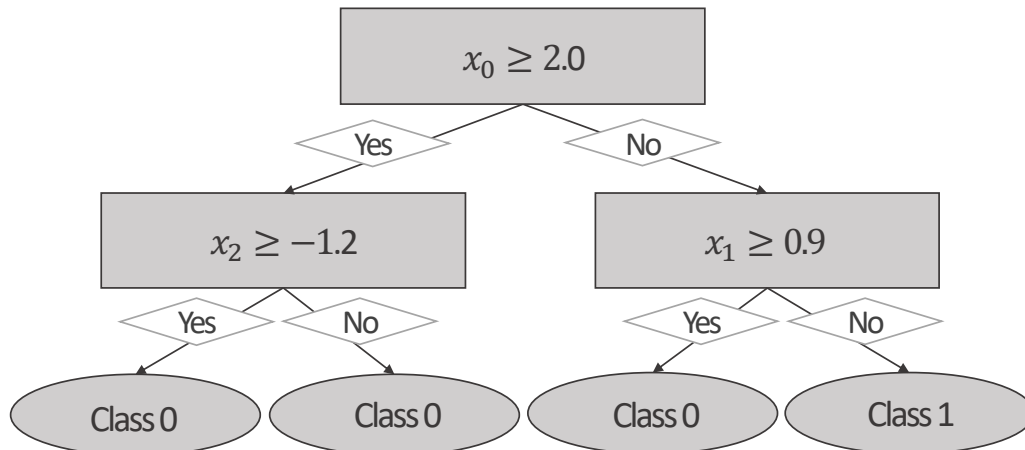
1. $S_0 * S_1 * \lambda_1$
2. $S_0 * (1 - S_1) * \lambda_2$
3. $(1 - S_0) * S_2 * \lambda_3$
4. $(1 - S_0) * (1 - S_2) * \lambda_4$

$$S_{\text{Heaviside}}(\mathbf{x}|\iota, \tau) = \begin{cases} 1, & \text{if } x_\iota \geq \tau \\ 0, & \text{otherwise} \end{cases}$$

GradTree: Gradient-Based Decision Trees

Dense Decision Tree Representation

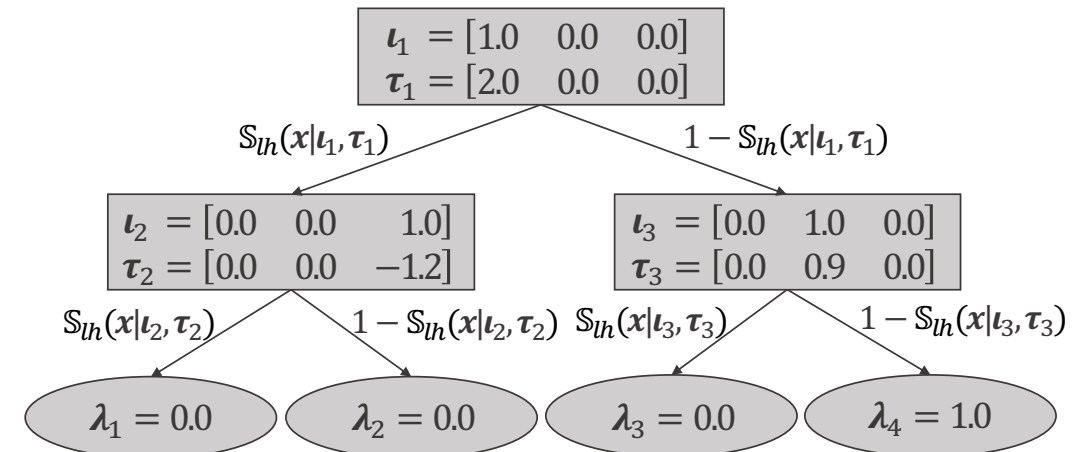
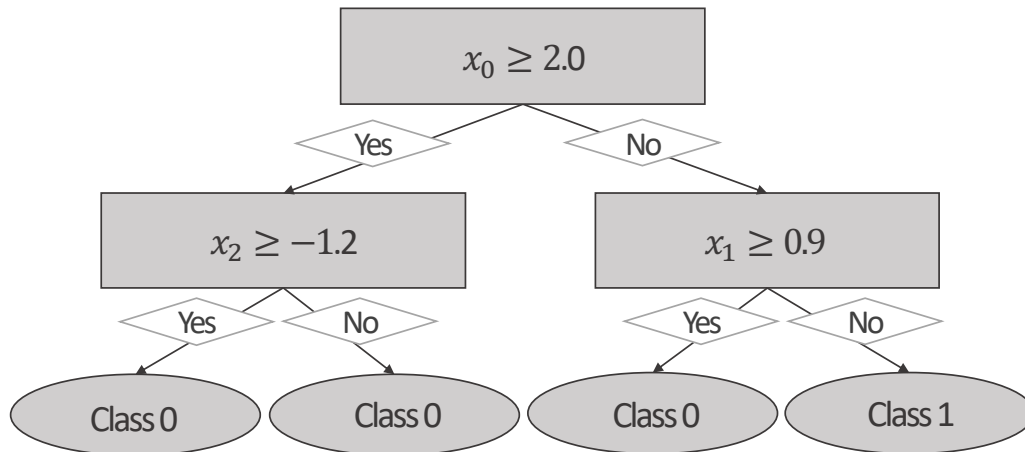
We propose a dense representation relaxing the split indices and split thresholds to account for the fact that feature indices are categorical instead of ordinal.



GradTree: Gradient-Based Decision Trees

Dense Decision Tree Representation

We propose a dense representation relaxing the split indices and split thresholds to account for the fact that feature indices are categorical instead of ordinal.

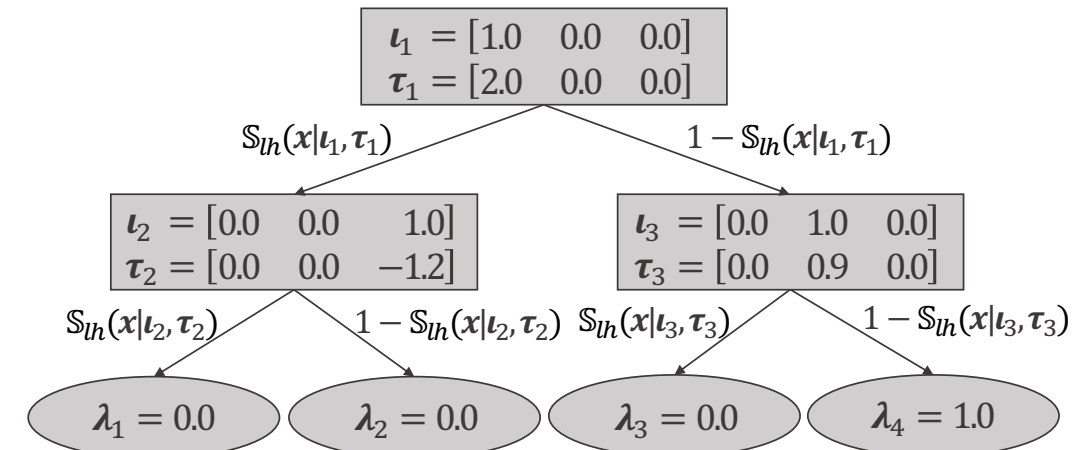
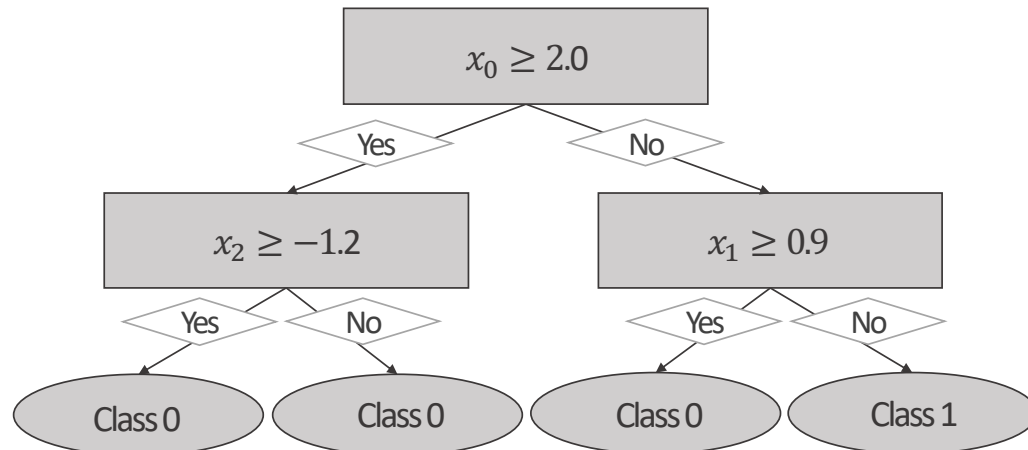


$$\begin{aligned}
 S_{\text{logistic}}(\mathbf{x}|\boldsymbol{\nu}, \boldsymbol{\tau}) &= \sigma \left(\sum_{i=0}^n \nu_i x_i - \sum_{i=0}^n \nu_i \tau_i \right) \\
 &= \sigma (\boldsymbol{\nu} \cdot \mathbf{x} - \boldsymbol{\nu} \cdot \boldsymbol{\tau})
 \end{aligned}$$

GradTree: Gradient-Based Decision Trees

Dense Decision Tree Representation

We propose a dense representation relaxing the split indices and split thresholds to account for the fact that feature indices are categorical instead of ordinal.



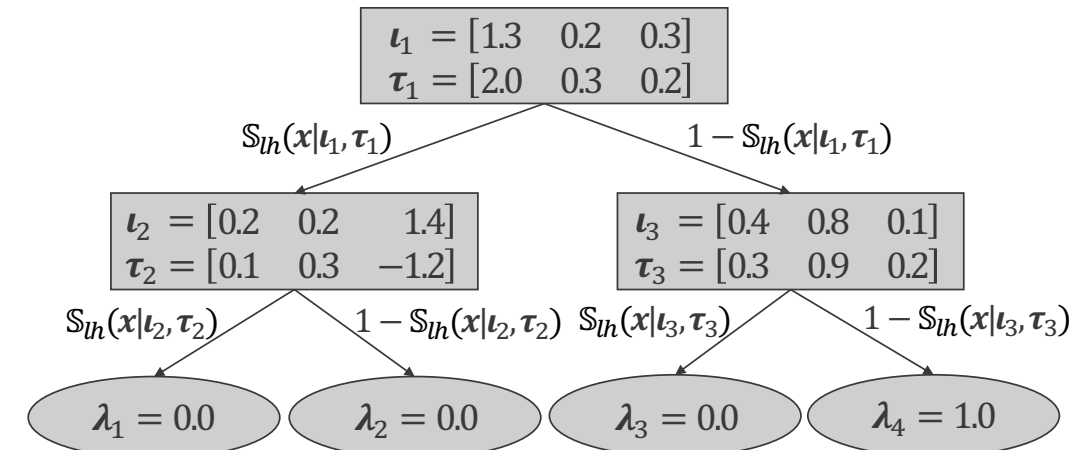
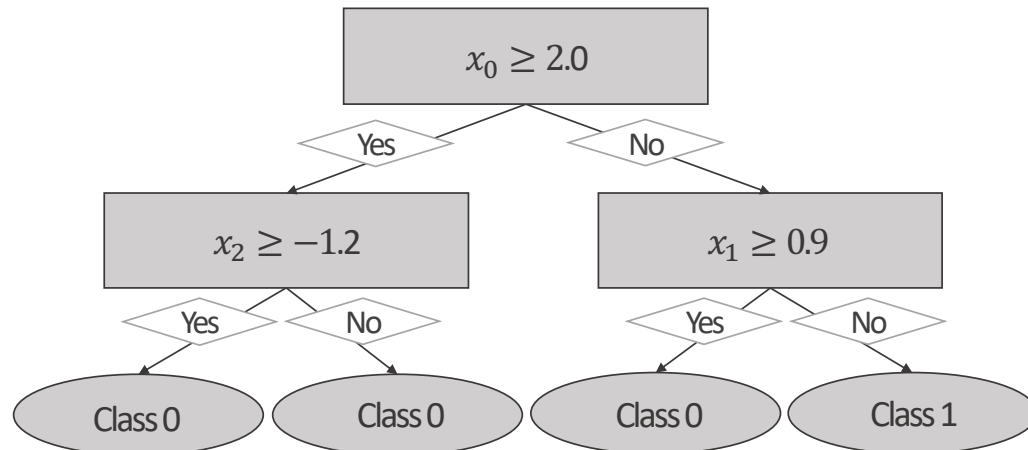
1. Round output of sigmoid function (utilized as split function) → **hard splits**

$$\begin{aligned}
 S_{\text{logistic}}(\mathbf{x}|\boldsymbol{\nu}, \boldsymbol{\tau}) &= \sigma \left(\sum_{i=0}^n \nu_i x_i - \sum_{i=0}^n \nu_i \tau_i \right) \\
 &= \sigma (\boldsymbol{\nu} \cdot \mathbf{x} - \boldsymbol{\nu} \cdot \boldsymbol{\tau})
 \end{aligned}$$

GradTree: Gradient-Based Decision Trees

Dense Decision Tree Representation

We propose a dense representation relaxing the split indices and split thresholds to account for the fact that feature indices are categorical instead of ordinal.



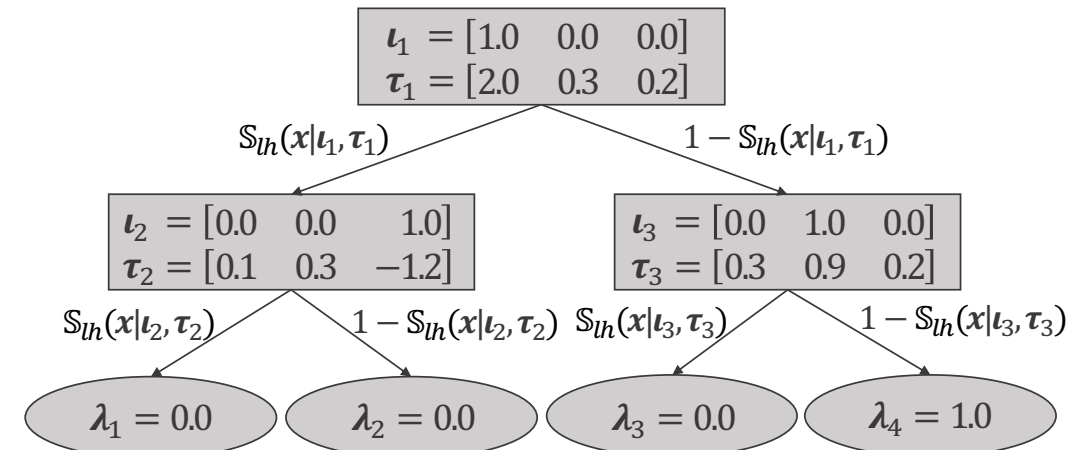
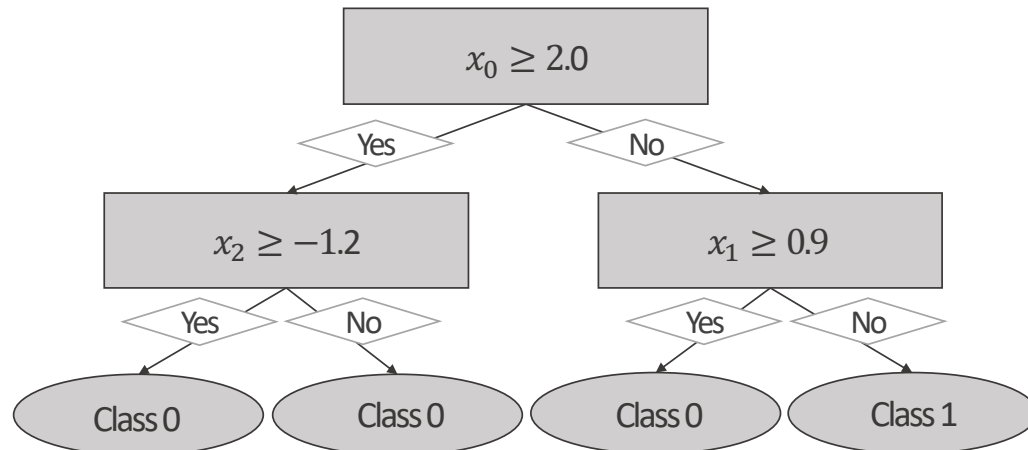
1. Round output of sigmoid function (utilized as split function) \rightarrow **hard splits**

$$\begin{aligned}
 S_{\text{logistic}}(\mathbf{x}|\boldsymbol{\nu}, \boldsymbol{\tau}) &= \sigma \left(\sum_{i=0}^n \nu_i x_i - \sum_{i=0}^n \nu_i \tau_i \right) \\
 &= \sigma(\boldsymbol{\nu} \cdot \mathbf{x} - \boldsymbol{\nu} \cdot \boldsymbol{\tau})
 \end{aligned}$$

GradTree: Gradient-Based Decision Trees

Dense Decision Tree Representation

We propose a dense representation relaxing the split indices and split thresholds to account for the fact that feature indices are categorical instead of ordinal.



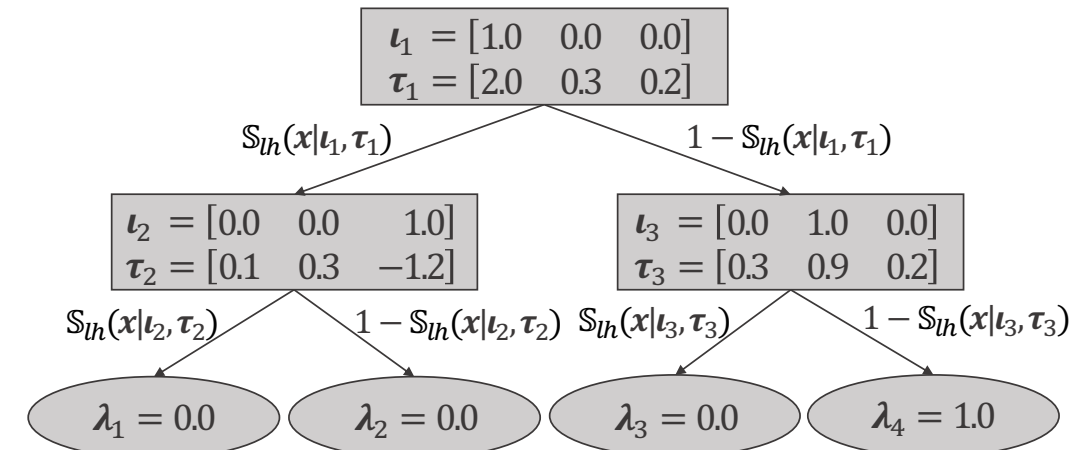
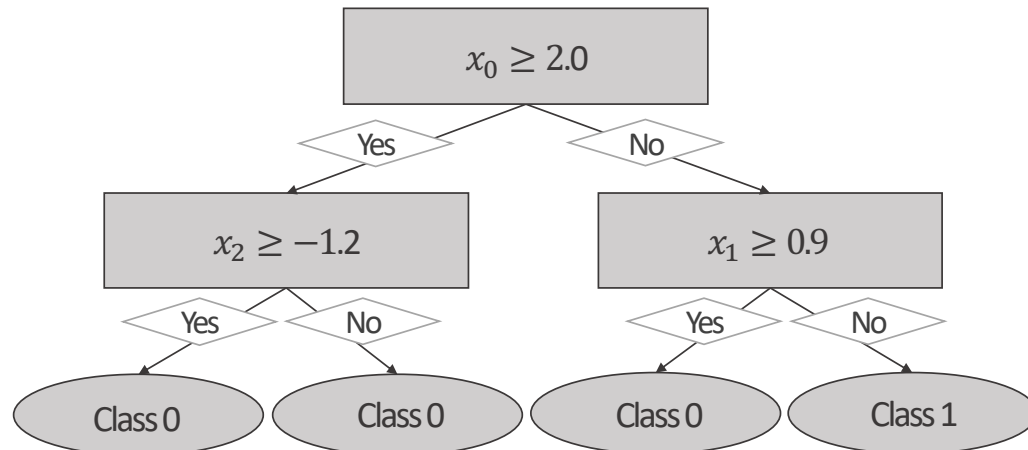
1. Round output of sigmoid function (utilized as split function) \rightarrow **hard splits**
2. Hardmax function to enforce one-hot encoded split index vectors \rightarrow **univariate DTs**

$$\begin{aligned}
 S_{\text{logistic}}(\mathbf{x}|\boldsymbol{\iota}, \boldsymbol{\tau}) &= \sigma \left(\sum_{i=0}^n \iota_i x_i - \sum_{i=0}^n \iota_i \tau_i \right) \\
 &= \sigma(\boldsymbol{\iota} \cdot \mathbf{x} - \boldsymbol{\iota} \cdot \boldsymbol{\tau})
 \end{aligned}$$

GradTree: Gradient-Based Decision Trees

Dense Decision Tree Representation

We propose a dense representation relaxing the split indices and split thresholds to account for the fact that feature indices are categorical instead of ordinal.



1. Round output of sigmoid function (utilized as split function) → **hard splits**
2. Hardmax function to enforce one-hot encoded split index vectors → **univariate DTs**

→ **Straight-Through (ST) Operator for Backpropagation of DT Loss to deal with non-differentiable operations**

$$\begin{aligned} S_{\text{logistic}}(\mathbf{x}|\boldsymbol{\nu}, \boldsymbol{\tau}) &= \sigma \left(\sum_{i=0}^n \nu_i x_i - \sum_{i=0}^n \nu_i \tau_i \right) \\ &= \sigma(\boldsymbol{\nu} \cdot \mathbf{x} - \boldsymbol{\nu} \cdot \boldsymbol{\tau}) \end{aligned}$$

Evaluation

- Best performance on most datasets
- Significantly outperformed gradient-based counterpart (DNDTs)
- Substantial performance increase on several datasets (e.g. *Echocardiogram & Absenteeism*)
- Good performance with default hyperparameters
- Robust to overfitting
- Small effective tree size
- Efficient for large and high-dimensional datasets (average runtime of 35 seconds)

Table 1: **Binary Classification Performance.** We report macro F1-scores (mean \pm stdev over 10 trials) on test data with optimized hyperparameters. The rank of each method is presented in brackets.

	Gradient-Based		Non-Greedy		Greedy
	GradTree (ours)	DNDT	GeneticTree	DL8.5 (Optimal)	CART
Blood Transfusion	0.628 \pm .036 (1)	0.543 \pm .051 (5)	0.575 \pm .094 (4)	0.590 \pm .034 (3)	0.613 \pm .044 (2)
Banknote Authentication	0.987 \pm .007 (1)	0.888 \pm .013 (5)	0.922 \pm .021 (4)	0.962 \pm .011 (3)	0.982 \pm .007 (2)
Titanic	0.776 \pm .025 (1)	0.726 \pm .049 (5)	0.730 \pm .074 (4)	0.754 \pm .031 (2)	0.738 \pm .057 (3)
Raisins	0.840 \pm .022 (4)	0.821 \pm .033 (5)	0.857 \pm .021 (1)	0.849 \pm .027 (3)	0.852 \pm .017 (2)
Rice	0.926 \pm .007 (3)	0.919 \pm .012 (5)	0.927 \pm .005 (2)	0.925 \pm .008 (4)	0.927 \pm .006 (1)
Echocardiogram	0.658 \pm .113 (1)	0.622 \pm .114 (3)	0.628 \pm .105 (2)	0.609 \pm .112 (4)	0.555 \pm .111 (5)
Wisconsin Breast Cancer	0.904 \pm .022 (2)	0.913 \pm .032 (1)	0.892 \pm .028 (4)	0.896 \pm .021 (3)	0.886 \pm .025 (5)
Loan House	0.714 \pm .041 (1)	0.694 \pm .036 (2)	0.451 \pm .086 (5)	0.607 \pm .045 (4)	0.662 \pm .034 (3)
Heart Failure	0.750 \pm .070 (3)	0.754 \pm .062 (2)	0.748 \pm .068 (4)	0.692 \pm .062 (5)	0.775 \pm .054 (1)
Heart Disease	0.779 \pm .047 (1)	$n > 12$	0.704 \pm .059 (4)	0.722 \pm .065 (2)	0.715 \pm .062 (3)
Adult	0.743 \pm .034 (2)	$n > 12$	0.464 \pm .055 (4)	0.723 \pm .011 (3)	0.771 \pm .011 (1)
Bank Marketing	0.640 \pm .027 (1)	$n > 12$	0.473 \pm .002 (4)	0.502 \pm .011 (3)	0.608 \pm .018 (2)
Congressional Voting	0.950 \pm .021 (1)	$n > 12$	0.942 \pm .021 (2)	0.924 \pm .043 (4)	0.933 \pm .032 (3)
Absenteeism	0.626 \pm .047 (1)	$n > 12$	0.432 \pm .073 (4)	0.587 \pm .047 (2)	0.564 \pm .042 (3)
Hepatitis	0.608 \pm .078 (2)	$n > 12$	0.446 \pm .024 (4)	0.586 \pm .083 (3)	0.622 \pm .078 (1)
German	0.592 \pm .068 (1)	$n > 12$	0.412 \pm .006 (4)	0.556 \pm .035 (3)	0.589 \pm .065 (2)
Mushroom	1.000 \pm .001 (1)	$n > 12$	0.984 \pm .003 (4)	0.999 \pm .001 (2)	0.999 \pm .001 (3)
Credit Card	0.674 \pm .014 (4)	$n > 12$	0.685 \pm .004 (1)	0.679 \pm .007 (3)	0.683 \pm .010 (2)
Horse Colic	0.842 \pm .039 (1)	$n > 12$	0.496 \pm .169 (4)	0.708 \pm .038 (3)	0.786 \pm .062 (2)
Thyroid	0.905 \pm .010 (2)	$n > 12$	0.605 \pm .116 (4)	0.682 \pm .018 (3)	0.922 \pm .011 (1)
Cervical Cancer	0.521 \pm .043 (1)	$n > 12$	0.514 \pm .034 (2)	0.488 \pm .027 (4)	0.506 \pm .034 (3)
Spambase	0.903 \pm .025 (2)	$n > 12$	0.863 \pm .019 (3)	0.863 \pm .011 (4)	0.917 \pm .011 (1)
Mean Relative Diff. (MRD) \downarrow	0.008 \pm .012 (1)	0.056 \pm .051 (3)	0.211 \pm .246 (5)	0.084 \pm .090 (4)	0.035 \pm .048 (2)
Mean Reciprocal Rank (MRR) \uparrow	0.758 \pm .306 (1)	0.370 \pm .268 (3)	0.365 \pm .228 (4)	0.335 \pm .090 (5)	0.556 \pm .293 (2)

Implementation

```
from GradTree import GradTree

params = {
    'depth': 5,

    'learning_rate_index': 0.01,
    'learning_rate_values': 0.01,
    'learning_rate_leaf': 0.005,

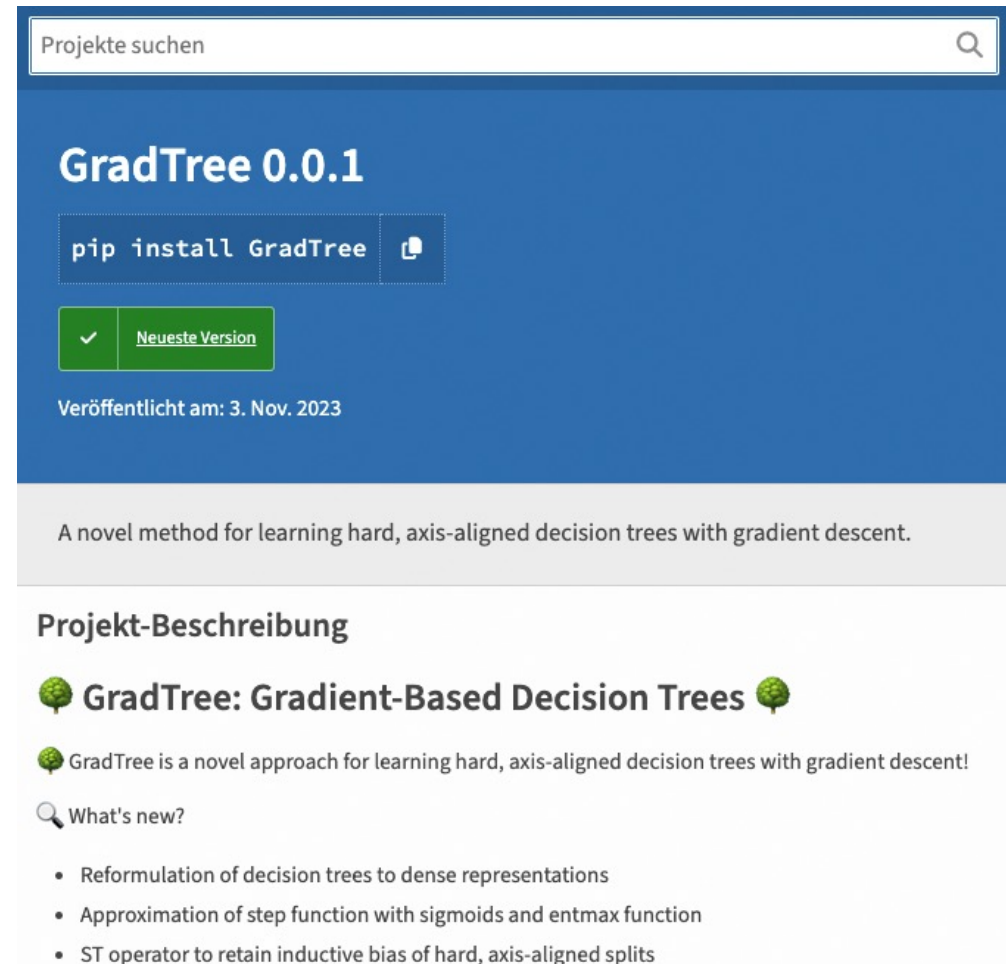
    'loss': 'crossentropy',
}

args = {
    'cat_idx': categorical_feature_indices,
    'objective': 'binary',
}

model_gradtree = GradTree(params=params, args=args)

model_gradtree.fit(X_train=X_train,
                  y_train=y_train,
                  X_val=X_valid,
                  y_val=y_valid)

model_gradtree = model_gradtree.predict(X_test)
```



Projekte suchen

GradTree 0.0.1



pip install GradTree


✓ Neueste Version


Veröffentlicht am: 3. Nov. 2023

A novel method for learning hard, axis-aligned decision trees with gradient descent.

Projekt-Beschreibung

 **GradTree: Gradient-Based Decision Trees** 

 GradTree is a novel approach for learning hard, axis-aligned decision trees with gradient descent!

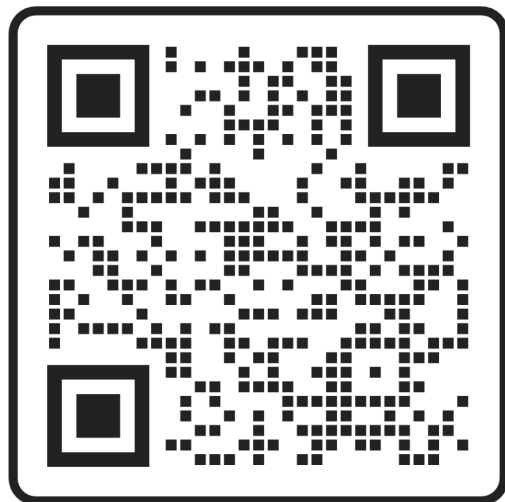
 What's new?

- Reformulation of decision trees to dense representations
- Approximation of step function with sigmoids and entmax function
- ST operator to retain inductive bias of hard, axis-aligned splits

<https://pypi.org/project/GRANDE/>

Thank you for your attention!

Paper:



Visit our poster (#2636) on
Saturday, February 24



Follow-Up Work @ICLR'24



**GRANDE: Gradient-Based Decision
Tree Ensembles**