# Feature learning for interpretable, Performant Decision Trees

January 8, 2024

Reviewr : Park Seok Hun

## Table of Contents

## Table of Contents

- Decision Tree :

    depth $\downarrow$    ►    Performance $\downarrow$, Interpretable $\uparrow$

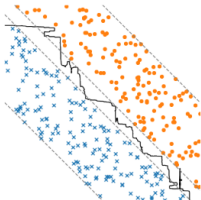    depth $\uparrow$    ►    Performance $\uparrow$, Interpretable $\downarrow$

- No matter how much the depth is increased, the performance on test data does not significantly outperform other models.

# Introduction

- They propose an algorithm that, through **feature learning**, generates a single tree with an appropriate depth for a Decision Tree while achieving good performance.

- **Feature learning** means that in the decision tree training, instead of $X_j \leq c$, the algorithm contemplates $f(\mathbf{X}) \leq c$, and learning $f$ during the training process where $\mathbf{X} = (X_1, ..., X_p)'$.
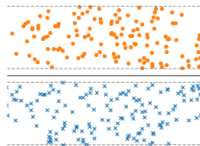
(a) A decision tree is relatively complex and generalizes poorly.

(b) A random forest is very complex and generalizes better, but not perfectly.

(c) After rotation, a decision tree is simple and generalizes perfectly.

- Above picture demonstrates that through feature transformation, a single tree can achieve good performance.

## Proposed method

- They propose alternating between learning a tree, similar to the CART, and performing feature learning based on gradients.

- For the feature learning, they consider Kernel Density Decision Tree(KDDT) which is differentiable.

## Table of Contents

## Fuzzy Decision Tree(FDT)

Problem of conventional Decision Tree(CART) :

Rule : If $X <= 3000$, $X$ is classified as group A else B.

Then, $X = 2999$ and $X = 3001$ are classified as different group.

- This causes prediction errors for Decision Trees near the boundaries
- If CART deterministically split data into child node, FDT reflect the possibility of data being split into each child node (e.g., using probabilities)

## Crisp Decision Tree

- Crisp Decision Tree : CART

## Table of Contents

## Kernel Density Decision Tree(KDDT)

- KDDT is a model that expresses the likelihood of splitting into child nodes using a kernel to represent probabilities.
- Note that unlike CART, KDDT is the differentiable model.

Example

Let $X \in \mathbb{R}$ be a input vector. Then, we have

$\mathbb{I}(X \in [a_j, b_j]) \to F(X - a_j) - F(X - b_j)$ where $F$ is cdf of normal.

## Table of Contents

- Any differentiable parameterized class of feature transforms can be used.
- example : Linear transformation $\mathbf{X} \rightarrow A\mathbf{X} + b$
- When training rule in the KDDT, feature learning is conducted based on gradient method for the impurity measure.

## Table of Contents

## Fuzzy into Crisp

- In the KDDT paper, they proposed the method for converting fuzzy decision trees to crisp decision trees, and it seems to be employed here.

- However, the details are not elaborated upon, and there is no code available.

- Note that the performance slightly degrades during the converting.
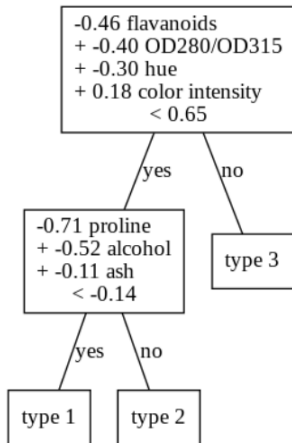
## Table of Contents

| data $n, p, q$ | LR | MLP | DT | RF | ET | XGB | ours: linear fuzzy | crisp |
|---|---|---|---|---|---|---|---|---|
| iris [18] | **0.960** | 0.953 | 0.947 | 0.947 | 0.953 | 0.947 | **0.960** | **0.960** |
| 150, 4 (4), 3 | - | - | 6.4 | 7.2e2 | 2.1e3 | 4.3e2 | 6.1 | 7.6 |
| heart-disease [30] | **0.822** | 0.792 | 0.707 | 0.802 | 0.795 | 0.792 | 0.812 | 0.812 |
| 303, 13 (20), 2 | - | - | 13.9 | 4.8e3 | 1.1e4 | 7.9e2 | 21.6 | 19.4 |
| dry-bean [31] | 0.925 | **0.934** | 0.912 | 0.923 | 0.921 | 0.928 | 0.920 | 0.913 |
| 13611, 16 (16), 7 | - | - | 99.8 | 6.7e4 | 2.0e5 | 1.3e4 | 1.1e2 | 45.8 |
| wine [1] | 0.983 | **0.989** | 0.904 | 0.977 | **0.989** | 0.955 | 0.983 | 0.983 |
| 178, 13 (13), 3 | - | - | 8.5 | 9.4e2 | 3.3e3 | 2.4e2 | 2.0 | 2.0 |
| car [5] | 0.926 | 0.992 | 0.977 | 0.964 | 0.971 | **0.994** | 0.991 | 0.992 |
| 1728, 6 (21), 4 | - | - | 95.3 | 2.3e4 | 3.1e4 | 4.5e3 | 29.0 | 29.0 |
| wdbc [44] | 0.974 | **0.975** | 0.935 | 0.965 | 0.970 | 0.968 | 0.972 | 0.972 |
| 569, 30 (30), 2 | - | - | 13.0 | 1.9e3 | 6.0e3 | 2.7e2 | 1.3 | 1.3 |
| sonar [38] | 0.755 | 0.879 | 0.735 | 0.826 | **0.880** | 0.855 | 0.818 | 0.799 |
| 208, 60 (60), 2 | - | - | 14.1 | 2.0e3 | 5.6e3 | 3.0e2 | 5.7 | 3.9 |
| pendigits [2] | 0.952 | **0.994** | 0.964 | 0.993 | **0.994** | 0.991 | 0.981 | 0.976 |
| 10992, 16 (16), 10 | - | - | 3.2e2 | 3.8e4 | 9.8e4 | 8.5e3 | 2.6e2 | 2.4e2 |
| ionosphere [39] | 0.875 | 0.917 | 0.892 | 0.934 | **0.943** | **0.943** | 0.932 | 0.920 |
| 351, 34 (34), 2 | - | - | 15.5 | 2.2e3 | 5.9e3 | 3.4e2 | 3.9 | 5.5 |

- Even as a single tree, proposed model performs well.

- The above results is from converted crisp proposed model.