

# SAT-based Approach for Learning Optimal Decision Trees with Non-Binary Features

---

POUYA SHATI, ELDAN COHEN, SHEILA MCILRAITH

UNIVERSITY OF TORONTO

CP 2021

SEPTEMBER 13, 2021

# Overview

---

- **Decision trees** are popular classification models
  - provide **interpretability** and **accuracy**
  - constructed via **greedy heuristics** or **exact methods**
  - exact optimization methods largely focus on **binary features**
- **Our contribution:** an approach to handle **non-binary** features effectively
  - outperforms the state of the art on **non-binary** datasets with two popular objectives

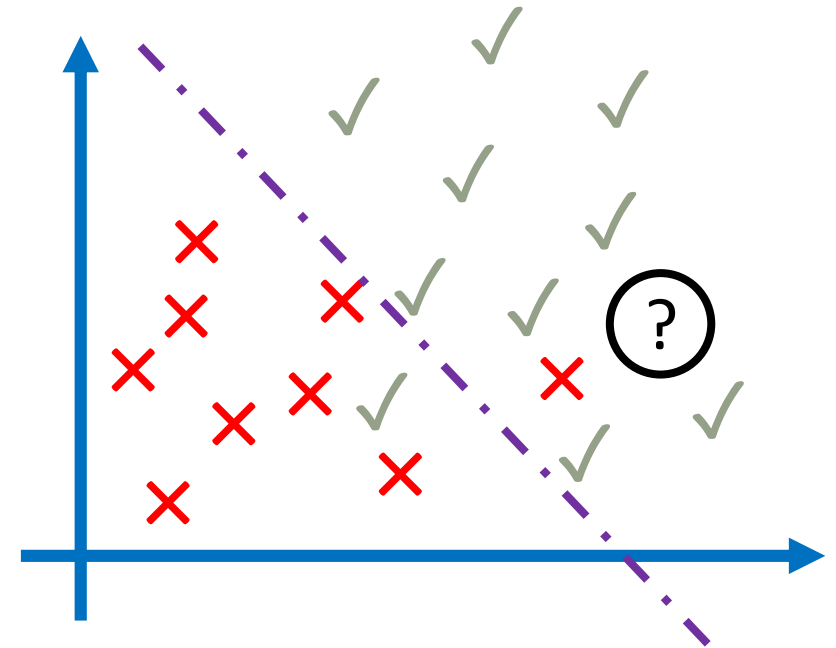
# Background

---

# Classification

---

- A popular application of **machine learning**
- Labelling function learned from labelled data set
- The goal is to achieve high accuracy on unseen data points



# Decision Trees

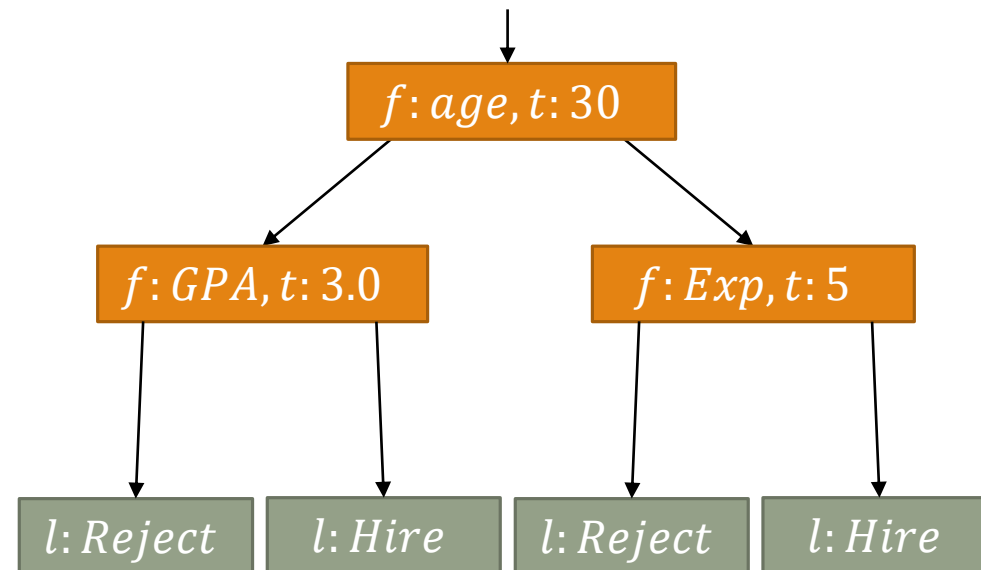
---

- **Decision trees are interpretable:**
  - human-readable
  - amenable to further (logical) reasoning
- Prime candidates for **safety-critical** applications

# Decision Trees

---

- **Branching nodes** perform a **split** on a given feature and threshold
- **Leaf nodes** assign a label



# Decision Trees

---

- A set of features  $F$  and integer labels  $C$
- A decision tree:  $\mathcal{D} = (\mathcal{T}, \beta, \alpha, \theta)$ :
  - $\mathcal{T}$  tree structure  $(\mathcal{T}_\beta, \mathcal{T}_L, \delta, p, l, r)$
  - $\beta$  feature selection function
  - $\alpha$  threshold selection function
  - $\theta$  leaf labelling function

- Recursive prediction for point  $x_i$ :

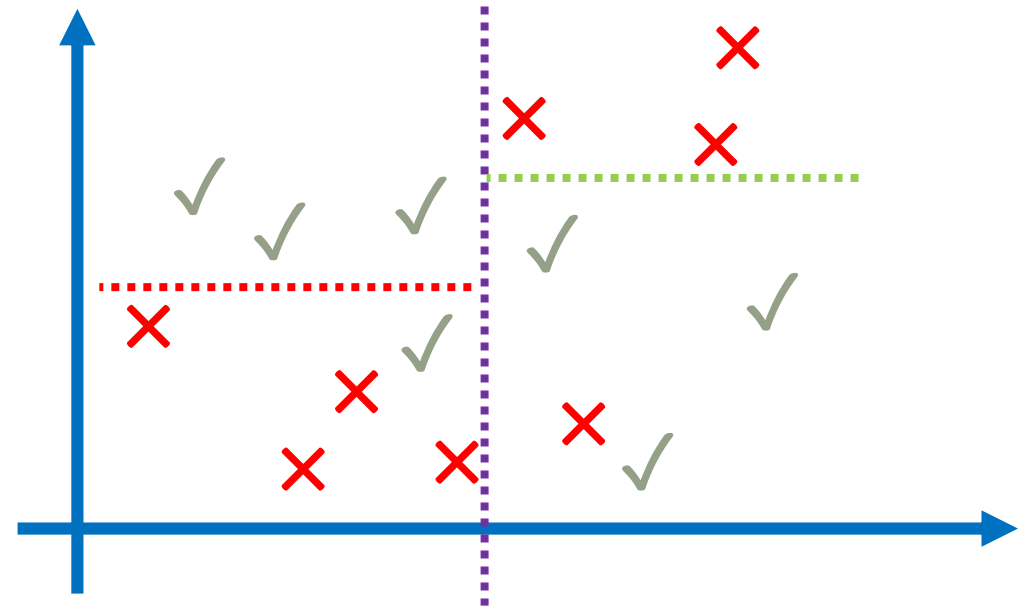
- $$\Theta(t, x_i) = \begin{cases} \theta(t) & \text{if } t \in \mathcal{T}_L \\ \Theta(l(t), x_i) & \text{else if } x_i[\beta(t)] \leq \alpha(t) \\ \Theta(r(t), x_i) & \text{else} \end{cases}$$

# Decision Trees

---

Ways to construct decision trees:

1. **Local search and heuristics**
2. **Combinatorial optimization:**
  - optimality guarantees
  - additional constraints





# Optimization Problem

---

## SAT:

- A set of variables  $\mathcal{X} = \{x_0, x_1, \dots, x_n\}$  and a set of clauses  $\mathcal{C} = \{C_1, C_2, \dots, C_k\}$
- Find an assignment  $\mathcal{M}: \mathcal{X} \rightarrow \{false, true\}$  that satisfies **all** clauses

## MaxSAT:

- All **hard** clauses  $\mathcal{C}_h$  should be satisfied
- The number of satisfied **soft** clauses  $\mathcal{C}_s$  needs to be **maximized**

# Encoding

---

# Encoding Components

---

- It is straight-forward to encode:
  - feature selection
  - leaf labelling
  - presence at leaves
- The challenging component is the **split**
- How can we model a numerical **threshold**?

# Split Encoding

---

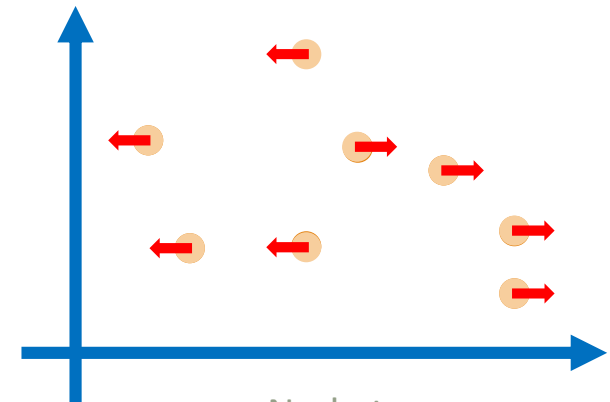
- Existing approach:
  - only support **binary** features!
  - transform **numerical** and **categorical** features into a set of **binary** ones
  - can lead to a huge number of features
- Avellaneda's [2020], Hu et al.'s [2020], and Verhaeghe et al.'s [2020] employ this approach

Numerical	Binary
4	001
9	011
1	000
4	001
12	111

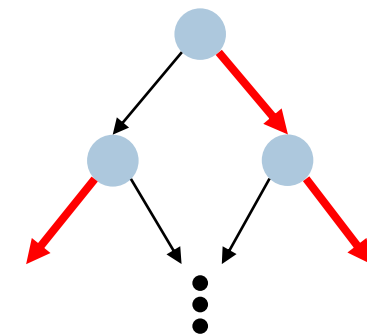
Categorical	Binary
A	001
B	010
A	001
C	100

# Split Encoding

- New idea:
  - encode the **direction** for each Point instead
  - validate the directions according to the **order** of values
  - the directions for **absent** points are encoded as well



Node t



Point i

# Our Encoding

---

## Variables:

- $a_{t,j}$ : feature  $j$  is chosen at node  $t$
- $s_{i,t}$ : point  $i$  is directed left at node  $t$
- $z_{i,t}$ : point  $i$  ends up at leaf  $t$
- $g_{t,c}$ : label  $c$  is assigned to leaf  $t$

## Parameters:

- set of features  $F$  and integer labels  $\mathcal{C}$
- set of training examples  $\mathcal{X}$
- labelling  $\gamma: X \rightarrow \mathcal{C}$
- tree structure  $T = \{\delta, T_B, T_L, p, l, r\}$

# Our Encoding

---

## Variables:

- $a_{t,j}$ : feature  $j$  is chosen at node  $t$
- $s_{i,t}$ : point  $i$  is directed left at node  $t$
- $z_{i,t}$ : point  $i$  ends up at leaf  $t$
- $g_{t,c}$ : label  $c$  is assigned to leaf  $t$

## Clauses:

- Exactly **one feature** is chosen at each **branching node**

$$\left( \neg a_{t,j}, \neg a_{t,j'} \right) \quad t \in \mathcal{T}_B, j \neq j' \in F$$

$$\left( \bigvee_{j \in F} a_{t,j} \right) \quad t \in \mathcal{T}_B$$

# Our Encoding

---

## Variables:

- $a_{t,j}$ : feature  $j$  is chosen at node  $t$
- $s_{i,t}$ : point  $i$  is directed left at node  $t$
- $z_{i,t}$ : point  $i$  ends up at leaf  $t$
- $g_{t,c}$ : label  $c$  is assigned to leaf  $t$

## Clauses:

- The **directions** for splits are in **order**

$$(\neg a_{t,j}, s_{i,t}, \neg s_{i',t}) \quad t \in \mathcal{T}_B, j \in F, (i, i') \in O_j(\mathcal{X})$$

$$(\neg a_{t,j}, \neg s_{i,t}, s_{i',t}) \quad t \in \mathcal{T}_B, j \in F, (i, i') \in O_j(\mathcal{X}), x_i[j] = x_{i'}[j]$$



# Our Encoding

---

## Variables:

- $a_{t,j}$ : feature  $j$  is chosen at node  $t$
- $s_{i,t}$ : point  $i$  is directed left at node  $t$
- $z_{i,t}$ : point  $i$  ends up at leaf  $t$
- $g_{t,c}$ : label  $c$  is assigned to leaf  $t$

## Clauses:

- The splits are **non-trivial**

$$\left( \neg a_{t,j}, s_{\#_j^1, t} \right) \quad t \in \mathcal{T}_B, j \in F$$

$$\left( \neg a_{t,j}, s_{\#_j^{|\mathcal{X}|}, t} \right) \quad t \in \mathcal{T}_B, j \in F$$

# Our Encoding

---

Variables:

- $a_{t,j}$ : feature  $j$  is chosen at node  $t$
- $s_{i,t}$ : point  $i$  is directed left at node  $t$
- $z_{i,t}$ : point  $i$  ends up at leaf  $t$
- $g_{t,c}$ : label  $c$  is assigned to leaf  $t$

Clauses:

- **Presence at leaves** matches the split directions

$$\left( \neg z_{i,t}, s_{i,t'} \right) \quad t \in \mathcal{T}_L, x_i \in \mathcal{X}, t' \in A_l(t)$$

$$\left( \neg z_{i,t}, \neg s_{i,t'} \right) \quad t \in \mathcal{T}_L, x_i \in \mathcal{X}, t' \in A_r(t)$$

$$\left( z_{i,t}, \bigvee_{t' \in A_l(t)} \neg s_{i,t'}, \bigvee_{t' \in A_r(t)} s_{i,t'} \right) \quad t \in \mathcal{T}_L, x_i \in \mathcal{X}$$

# Our Encoding

---

## Variables:

- $a_{t,j}$ : feature  $j$  is chosen at node  $t$
- $s_{i,t}$ : point  $i$  is directed left at node  $t$
- $z_{i,t}$ : point  $i$  ends up at leaf  $t$
- $g_{t,c}$ : label  $c$  is assigned to leaf  $t$

## Clauses:

- At most one **label** is chosen at each leaf

$$(\neg g_{t,c}, \neg g_{t,c'}) \quad t \in \mathcal{T}_L, c \neq c' \in \mathcal{C}$$

# Learning Decision Trees

---

Two main **objectives**:

- **Min-depth:**
  - correctly classify **all** of the training points
  - find the **lowest depth** possible
  - solved by iterative **SAT** instances
- **Max-accuracy:**
  - **maximize** the number of correct classifications
  - use a fixed **depth**
  - solved via **MaxSAT**

# Extension to Categorical Features

---

- Use the same idea for **categorical splits**:
  - no need to validate **order** in directions, checking **equality** is enough
  - enables **power set** branching:
    - **min-depth**: potentially more **shallow** solution
    - **max-accuracy**: potentially more **accurate** solution

# Experimental

---

# Experimental Setup

---

Objective	Language	Solver	Baseline	Baseline's Solver
Min-Depth	C++	MiniSAT	<b>Avellaneda's</b> [2020] SAT-based approach	MiniSAT
Max-Accuracy	Java	Loandra	<b>Hu et al.'s</b> [2020] MaxSAT approach	Loandra
			<b>Verhaeghe et al.'s</b> [2020] Constraint Programming approach	Oscar

- The chosen baselines are the state-of-the-art algorithms for their respective objectives

# Goals

---

- The benefits and applications of the two **objectives** are well-studied
- Focus on optimization performance:
  - find the solutions **faster**
  - find **near-optimal** solutions in **time-out** scenarios



# Datasets

---

- Three types of datasets:
  - mostly **numerical** features
  - mostly **categorical** features
  - mostly **binary** features

Type	Name	$ X $	$ F_N $	$ F_B $	$ F_C $	$\tilde{f}$	$ C $
N	Banknote	1372	4	0	0	5016	2
	Breast Cancer	116	9	0	0	891	2
	Cryotherapy	90	5	1	0	93	2
	Immunotherapy	91	6	1	0	166	2
	Ionosphere	351	32	2	0	8114	2
	Iris	150	4	0	0	119	3
	User Knowledge	258	5	0	0	431	4
	Vertebral Column	310	6	0	0	1741	2
	Wine	178	13	0	0	1263	3
B	Car <sup>†</sup>	1728	6	0	0	15	2
	Monk2	169	4	2	0	11	2

# Min-Depth Results

- On **non-binary** datasets, our approach is significantly faster than the baseline
- As expected, the existing approach works better on **binary** datasets

Dataset	Min Depth	Time (s)	
		Ours	SAT [3]
Banknote	4	<b>5.82</b>	T/O [4]
Breast Cancer	4	<b>6.59</b>	T/O [4]
Cryotherapy	4	<b>0.08</b>	0.24
Immunotherapy	4	<b>0.18</b>	1.3
Ionosphere	?	T/O [4]	T/O [3]
Iris	4	<b>0.04</b>	0.17
User Knowledge	5	<b>1.31</b>	59.44
Vertebral Column	5	<b>87.35</b>	T/O [5]
Wine	3	<b>0.11</b>	14.75
Car	8	T/O [8]	<b>89.1</b>
Monk2	6	2.73	<b>0.28</b>

# Max-Accuracy Results

- On **non-binary** datasets, our approach is significantly faster than the baselines
- As expected, existing approaches work better on **binary** datasets
- Our approach still finds **optimal solutions** for **binary** datasets most of the time

Dataset	Depth	Solution Cost			Time (s)		
		Ours	MaxSAT [15]	CP [23]	Ours	MaxSAT [15]	CP [23]
Banknote	2	<b>100</b>	176	<b>100</b>	<b>16.83</b>	T/O	512.21
	3	<b>23</b>	550	100	<b>105.79</b>	T/O	T/O
	4	<b>0</b>	88	100	<b>18.98</b>	T/O	T/O
Breast Cancer	2	<b>19</b>	24	<b>19</b>	<b>5.07</b>	T/O	22.19
	3	<b>9</b>	25	12	<b>242.16</b>	T/O	T/O
	4	<b>0</b>	18	11	<b>20.79</b>	T/O	T/O
Cryotherapy	2	<b>5</b>	<b>5</b>	<b>5</b>	<b>0.57</b>	3.68	4.21
	3	<b>1</b>	<b>1</b>	<b>1</b>	<b>0.73</b>	17.57	27.39
	4	<b>0</b>	<b>0</b>	<b>0</b>	<b>0.75</b>	24.14	7.61
Immunotherapy	2	<b>8</b>	<b>8</b>	<b>8</b>	<b>0.99</b>	10.53	5.22
	3	<b>4</b>	<b>4</b>	<b>4</b>	<b>3.81</b>	T/O	146.45
	4	<b>0</b>	<b>1</b>	<b>0</b>	<b>1.27</b>	T/O	18.53
Ionosphere	2	<b>29</b>	41	<b>29</b>	<b>155.06</b>	T/O	T/O
	3	<b>21</b>	186	29	T/O	T/O	T/O
	4	<b>10</b>	76	28	T/O	T/O	T/O
Iris	2	<b>6</b>	-	-	<b>0.6</b>	-	-
	3	<b>1</b>	-	-	<b>0.77</b>	-	-
	4	<b>0</b>	-	-	<b>0.82</b>	-	-
User Knowledge	2	<b>35</b>	-	-	<b>1.94</b>	-	-
	3	<b>10</b>	-	-	<b>3.29</b>	-	-
	4	<b>1</b>	-	-	<b>3.86</b>	-	-
Vertebral Column	2	<b>45</b>	46	<b>45</b>	<b>15.79</b>	T/O	67.91
	3	<b>32</b>	44	42	T/O	T/O	T/O
	4	<b>15</b>	39	42	T/O	T/O	T/O
Wine	2	<b>6</b>	-	-	<b>1.25</b>	-	-
	3	<b>0</b>	-	-	<b>1.62</b>	-	-
Car	2	<b>250</b>	<b>250</b>	<b>250</b>	12.67	9.2	<b>2.16</b>
	3	<b>182</b>	<b>182</b>	<b>182</b>	T/O	T/O	<b>5.99</b>
	4	<b>122</b>	<b>122</b>	<b>122</b>	T/O	T/O	<b>14.09</b>
Monk2	2	<b>57</b>	<b>57</b>	<b>57</b>	2.74	4.38	1.38
	3	<b>42</b>	<b>42</b>	<b>42</b>	T/O	826.31	<b>3.6</b>
	4	32	<b>31</b>	<b>31</b>	T/O	T/O	<b>8.12</b>

# Summary

---

- Novel MaxSAT-based encoding for constructing **optimal decision trees** for datasets with **numerical** and **categorical** Features
- Can be employed by both **min-depth** and **max-accuracy** objectives
- Supports **power set** splitting on **categorical** features to achieve **compactness**
- Significantly outperforms the state of the art for **non-binary** datasets

# Thank You!

## Questions & Answers

---

POUYA SHATI, ELDAN COHEN, SHEILA MCILRAITH

UNIVERSITY OF TORONTO

CP 2021

SEPTEMBER 13, 2021