

# Online Learning of Decision Trees with Thompson Sampling

Ayman Chaouki <sup>1 2</sup>   Jesse Read <sup>1</sup>   Albert Bifet <sup>2</sup>

<sup>1</sup>LIX, Ecole Polytechnique, IP Paris

<sup>2</sup>AI Institute, The University of Waikato

AISTATS 2022

# Table of Contents

- 1 Introduction
- 2 THOMPSON SAMPLING DECISION TREES (TSDT)
- 3 Experiments
- 4 Limitations and Future Work

# Section 1

## Introduction

# Decision Trees

- Hierarchical rules partitioning the feature space.
- Each element of the partition is assigned a class label.
- Interpretable models, useful for sensitive domains.

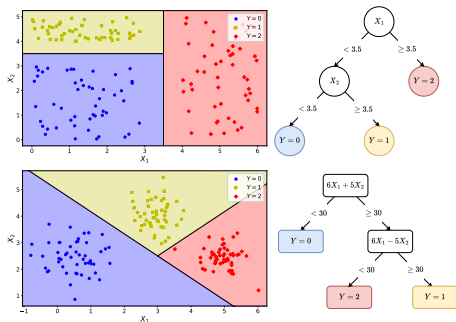


Figure: Example of Decision Tree classifiers on the feature space  $\mathbb{R}^2$ .

## Objective

Complex Decision Trees (i.e with many rules) lose interpretability. As such, it is desirable to find optimal Decision Trees both in the sense of maximising accuracy and also minimising complexity.

## Objective

Complex Decision Trees (i.e with many rules) lose interpretability. As such, it is desirable to find optimal Decision Trees both in the sense of maximising accuracy and also minimising complexity.

⚠ Unfortunately, this optimisation problem is NP-Hard.

- The continuous storage of data requires larger and larger memories to keep up.
- Many modern applications have an influx of data.
- Batch models are obsolete for handling data streams.

## Objective:

- Find optimal Decision Trees in the online setting.
- We formulate the problem within a Reinforcement Learning framework.
- We devise TSDT to solve the Reinforcement Learning problem.

## Section 2

# THOMPSON SAMPLING DECISION TREES (TSDT)



### Online Classification problem

- $X \in \mathcal{X} = \prod_{i=1}^q \{1, \dots, C_i\}$ ,  $Y \in \{1, \dots, K\}$ .
- Data  $(X_i, Y_i)$  arrive incrementally, and are i.i.d. following a joint probability distribution.

### Online Classification problem

- $X \in \mathcal{X} = \prod_{i=1}^q \{1, \dots, C_i\}$ ,  $Y \in \{1, \dots, K\}$ .
- Data  $(X_i, Y_i)$  arrive incrementally, and are i.i.d. following a joint probability distribution.
- For a Decision Tree  $T$  and  $X \in \mathcal{X}$ ,  $l(X)$  is the leaf containing  $X$ .
- Prediction:  $T(X) = \operatorname{Argmax}_{1 \leq k \leq K} \mathbb{P}[Y = k | X \in l(X)]$ .

### Online Classification problem

- $X \in \mathcal{X} = \prod_{i=1}^q \{1, \dots, C_i\}$ ,  $Y \in \{1, \dots, K\}$ .
- Data  $(X_i, Y_i)$  arrive incrementally, and are i.i.d. following a joint probability distribution.
- For a Decision Tree  $T$  and  $X \in \mathcal{X}$ ,  $l(X)$  is the leaf containing  $X$ .
- Prediction:  $T(X) = \operatorname{Argmax}_{1 \leq k \leq K} \mathbb{P}[Y = k | X \in l(X)]$ .
- Accuracy:  $\mathcal{H}(T) = \mathbb{P}[T(X) = Y]$ .

### Maximising Accuracy is not enough

The complete Decision Tree, that partitions  $\mathcal{X}$  into its unit-subsets, maximises Accuracy. Nonetheless, it is an uninterpretable solution that just corresponds to point-wise classification.

# Problem Formulation

## Objective function

### Regularised objective function

We seek to maximise Accuracy while penalising solutions with large complexities.

$$\mathcal{H}_r(T) = \mathcal{H}(T) - \lambda \mathcal{S}(T)$$

- $\lambda \geq 0$  is a penalty parameter.
- $\mathcal{S}(T)$  is the number of splits of  $T$ .
- Find  $T^*$  maximising  $\mathcal{H}_r(T)$ .

### Remark

This objective is used in CART during pruning. It is also employed by the OCT algorithm.

# Problem Formulation

## Markov Decision Process (MDP)

# Problem Formulation

## Markov Decision Process (MDP)

- **State:** Decision Tree.  $R$  denotes the root, the initial state.

# Problem Formulation

## Markov Decision Process (MDP)

- **State:** Decision Tree.  $R$  denotes the root, the initial state.
- **Action:** two types of actions: split actions split a leaf with respect to some attribute, and the terminal action ends the episode.

# Problem Formulation

## Markov Decision Process (MDP)

- **State:** Decision Tree.  $R$  denotes the root, the initial state.
- **Action:** two types of actions: split actions split a leaf with respect to some attribute, and the terminal action ends the episode.
- **Transition Dynamics:** When taking an action, we transition from state  $T$  to state  $T'$  deterministically, we denote the transition  $T \rightarrow T'$ . The set of next states from  $T$  is denoted  $\text{Ch}(T)$ . We also denote  $T \rightarrow \bar{T}$  when taking the terminal action.



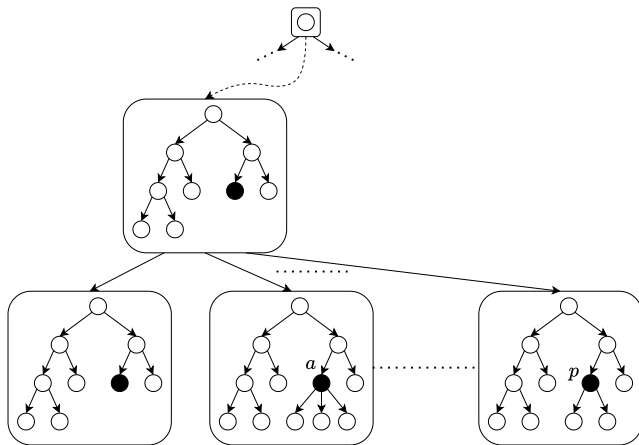
# Problem Formulation

## Markov Decision Process (MDP)

- **State:** Decision Tree.  $R$  denotes the root, the initial state.
- **Action:** two types of actions: split actions split a leaf with respect to some attribute, and the terminal action ends the episode.
- **Transition Dynamics:** When taking an action, we transition from state  $T$  to state  $T'$  deterministically, we denote the transition  $T \rightarrow T'$ . The set of next states from  $T$  is denoted  $\text{Ch}(T)$ . We also denote  $T \rightarrow \bar{T}$  when taking the terminal action.
- **Reward:** Split actions yield  $-\lambda$  reward, the terminal action yields reward  $\mathcal{H}(T)$ , which is unknown.

# Problem Formulation

## Search Tree representation



**Figure:** Graph representation of the state-action space. The (squared) nodes represent states and the edges represent actions.

# Problem Formulation

Linking the MDP to online Decision Tree optimisation

- Let  $\pi$  be a policy,  $\forall T : \pi(T) \in \text{Ch}(T)$ .
- $T = T^{(0)} \xrightarrow{\pi} T^{(1)} \xrightarrow{\pi} \dots, \xrightarrow{\pi} T^{(N)} = \overline{T^{(N-1)}}$ :

$$\mathcal{V}^\pi(T) = \mathbb{E} \left[ \sum_{j=1}^N r(T^{(j-1)}, T^{(j)}) \right]$$

## Proposition

The optimal policy  $\pi^* = \text{Argmax}_\pi \mathcal{V}^\pi(R)$  exists and it leads to  $T^*$ :

$$R \xrightarrow{\pi^*} \dots, \xrightarrow{\pi^*} T^* \xrightarrow{\pi^*} \overline{T^*}$$

$$T^* = \text{Argmax}_T \mathcal{H}_r(T)$$

## Bellman Optimality Equation

In any non-terminal state  $T$ , the optimal policy satisfies:

$$v^{\pi^*}(T) = \max_{T' \in \text{Ch}(T)} \left\{ -\lambda \mathbb{1}\{T' \neq \bar{T}\} + v^{\pi^*}(T') \right\}$$

- By estimating  $v^{\pi^*}(\bar{T})$ , TSDT estimates  $v^{\pi^*}(T)$  for non-terminal states by Backward Induction.

## Bellman Optimality Equation

In any non-terminal state  $T$ , the optimal policy satisfies:

$$\mathcal{V}^{\pi^*}(T) = \max_{T' \in \text{Ch}(T)} \left\{ -\lambda \mathbb{1}\{T' \neq \bar{T}\} + \mathcal{V}^{\pi^*}(T') \right\}$$

- By estimating  $\mathcal{V}^{\pi^*}(\bar{T})$ , TSDT estimates  $\mathcal{V}^{\pi^*}(T)$  for non-terminal states by Backward Induction.
- TSDT estimates  $\mathcal{V}^{\pi^*}(\bar{T})$  within a Bayesian framework, and backpropagate the posterior distributions up the Search Tree.
- TSDT's exploration policy is Thompson Sampling.

$$\mathcal{V}^{\pi^*}(\bar{T}) = \mathcal{H}(T) = \sum_{I \in \mathcal{L}(T)} \mathbb{P}[X \in I] \mathbb{E}[\mathbb{1}\{T(X) = Y\} | X \in I]$$

Posterior distribution on  $\mathcal{V}^{\pi^*}(\bar{T})$ :

$$\theta_{\bar{T}} = \sum_{I \in \mathcal{L}(T)} \hat{p}(I) \theta_{T,I}$$

$$\mathcal{V}^{\pi^*}(\bar{T}) = \mathcal{H}(T) = \sum_{I \in \mathcal{L}(T)} \mathbb{P}[X \in I] \mathbb{E}[\mathbb{1}\{T(X) = Y\} | X \in I]$$

Posterior distribution on  $\mathcal{V}^{\pi^*}(\bar{T})$ :

$$\theta_{\bar{T}} = \sum_{I \in \mathcal{L}(T)} \hat{p}(I) \theta_{T,I}$$

$$\theta_{T,I} \sim \text{Beta}(\alpha_{T,I}, \beta_{T,I})$$

$$\alpha_{T,I} = 1 + \sum_{i=1}^N \mathbb{1}\{X_i \in I\} \mathbb{1}\{T(X_i) = Y_i\};$$

$$\beta_{T,I} = 1 + \sum_{i=1}^N \mathbb{1}\{X_i \in I\} - \sum_{i=1}^N \mathbb{1}\{X_i \in I\} \mathbb{1}\{T(X_i) = Y_i\}$$

$$\mathcal{V}^{\pi^*}(\bar{T}) = \mathcal{H}(T) = \sum_{I \in \mathcal{L}(T)} \mathbb{P}[X \in I] \mathbb{E}[\mathbb{1}\{T(X) = Y\} | X \in I]$$

Posterior distribution on  $\mathcal{V}^{\pi^*}(\bar{T})$ :

$$\theta_{\bar{T}} = \sum_{I \in \mathcal{L}(T)} \hat{p}(I) \theta_{T,I}$$

$$\theta_{T,I} \sim \text{Beta}(\alpha_{T,I}, \beta_{T,I})$$

$$\alpha_{T,I} = 1 + \sum_{i=1}^N \mathbb{1}\{X_i \in I\} \underbrace{\mathbb{1}\{T(X_i) = Y_i\}}_{\text{unknown}};$$

$$\beta_{T,I} = 1 + \sum_{i=1}^N \mathbb{1}\{X_i \in I\} - \sum_{i=1}^N \mathbb{1}\{X_i \in I\} \underbrace{\mathbb{1}\{T(X_i) = Y_i\}}_{\text{unknown}}$$



$$T(X) = \text{Argmax}_{1 \leq k \leq K} \mathbb{P}[Y = k | X \in I(X)]$$

$$\forall I \in \mathcal{L}(T) : \hat{\rho}_k^{(i)}(I) = \frac{\sum_{j=1}^i \mathbb{1}\{X_j \in I\} \mathbb{1}\{Y_j = k\}}{\sum_{j=1}^i \mathbb{1}\{X_j \in I\}}$$

$$\hat{T}_i(I) = \text{Argmax}_{1 \leq k \leq K} \{\hat{\rho}_k^{(i)}(I)\}$$

$$T(X) = \text{Argmax}_{1 \leq k \leq K} \mathbb{P}[Y = k | X \in I(X)]$$

$$\forall I \in \mathcal{L}(T) : \hat{\rho}_k^{(i)}(I) = \frac{\sum_{j=1}^i \mathbb{1}\{X_j \in I\} \mathbb{1}\{Y_j = k\}}{\sum_{j=1}^i \mathbb{1}\{X_j \in I\}}$$

$$\hat{T}_i(I) = \text{Argmax}_{1 \leq k \leq K} \{\hat{\rho}_k^{(i)}(I)\}$$

$$\alpha_{T,I} = 1 + \sum_{i=2}^N \mathbb{1}\{X_i \in I\} \mathbb{1}\{\hat{T}_{i-1}(X_i) = Y_i\};$$

$$\beta_{T,I} = 1 + \sum_{i=2}^N \mathbb{1}\{X_i \in I\} - \sum_{i=2}^N \mathbb{1}\{X_i \in I\} \mathbb{1}\{\hat{T}_{i-1}(X_i) = Y_i\}$$

② What is the distribution of  $\theta_{\bar{T}}$ ?

② **What is the distribution of  $\theta_{\bar{T}}$ ?**

Posterior distribution on  $\mathcal{V}^{\pi^*}(\bar{T})$ :

$$\theta_{\bar{T}} = \sum_{I \in \mathcal{L}(T)} \hat{p}(I) \theta_{T,I}$$

$$\forall I \in \mathcal{L}(T) : \theta_{T,I} \sim \text{Beta}(\alpha_{T,I}, \beta_{T,I})$$

- $\theta_{\bar{T}}$  is a linear combination of independent Beta variables.

② **What is the distribution of  $\theta_{\bar{T}}$ ?**

Posterior distribution on  $\mathcal{V}^{\pi^*}(\bar{T})$ :

$$\theta_{\bar{T}} = \sum_{I \in \mathcal{L}(T)} \hat{p}(I) \theta_{T,I}$$

$$\forall I \in \mathcal{L}(T) : \theta_{T,I} \sim \text{Beta}(\alpha_{T,I}, \beta_{T,I})$$

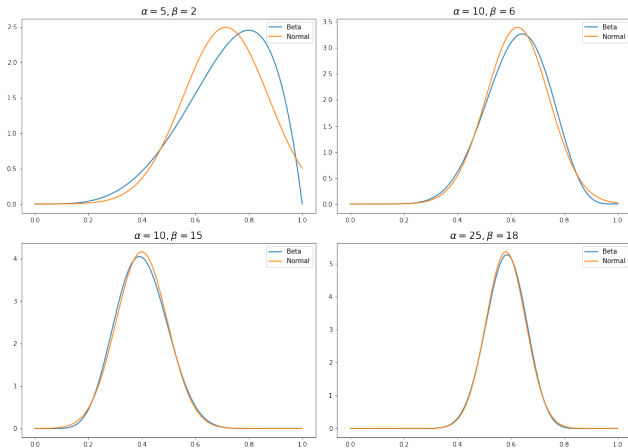
- $\theta_{\bar{T}}$  is a linear combination of independent Beta variables.
- Approximate each Beta distribution with a Normal distribution.

Match the first and second moments of the Beta distribution with its Normal approximation.

$$\theta_{T,l} \sim \mathcal{N} \left( \mu_{T,l}, (\sigma_{T,l})^2 \right)$$

$$\mu_{T,l} = \frac{\alpha_{T,l}}{\alpha_{T,l} + \beta_{T,l}}$$

$$(\sigma_{T,l})^2 = \frac{\alpha_{T,l}\beta_{T,l}}{(\alpha_{T,l} + \beta_{T,l})^2 (1 + \alpha_{T,l} + \beta_{T,l})}$$



**Figure:** Probability density functions of a Beta distribution and its Normal approximation for different parameters  $\alpha$ , and  $\beta$ .

$\theta_{\bar{T}} = \sum_{l \in \mathcal{L}(T)} \hat{p}(l) \theta_{T,l}$  linear combination of Normal variables:

$$\theta_{\bar{T}} \sim \mathcal{N} \left( \mu_{\bar{T}}, (\sigma_{\bar{T}})^2 \right)$$

$$\mu_{\bar{T}} = \sum_{l \in \mathcal{L}(T)} \hat{p}(l) \mu_{T,l}$$

$$(\sigma_{\bar{T}})^2 = \sum_{l \in \mathcal{L}(T)} \hat{p}(l)^2 (\sigma_{T,l})^2$$



## Bellman Optimality Equation

In any non-terminal state  $T$ , the optimal policy satisfies:

$$\mathcal{V}^{\pi^*}(T) = \max_{T' \in \text{Ch}(T)} \left\{ -\lambda \mathbb{1}\{T' \neq \bar{T}\} + \mathcal{V}^{\pi^*}(T') \right\}$$

$$\theta_T = \max_{T' \in \text{Ch}(T)} \left\{ -\lambda \mathbb{1}\{T' \neq \bar{T}\} + \theta_{T'} \right\}$$

$$\tilde{T} = \text{Argmax}_{T' \in \text{Ch}(T)} \{ -\lambda \mathbb{1}\{T' \neq \bar{T}\} + \mu_{T'} \}$$

$$\theta_T \sim \mathcal{N} \left( -\lambda \mathbb{1}\{\tilde{T} \neq \bar{T}\} + \mu_{\tilde{T}}, (\sigma_{\tilde{T}})^2 \right)$$

# TSDT

## The Algorithm

$$T^{(0)} = R \quad \square$$

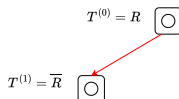
Initialise  $\pi(T) = \bar{T}$

**Iteration**  $t = 1$

- *Selection*
- *Simulation*
- *Expansion*
- *Backpropagation*

# TSDT

## The Algorithm



Initialise  $\pi(T) = \bar{T}$

**Iteration**  $t = 1$

- *Selection*
- *Simulation*
- *Expansion*
- *Backpropagation*

# TSDT

## The Algorithm



Initialise  $\pi(T) = \bar{T}$

**Iteration**  $t = 1$

- *Selection*
- *Simulation*
- *Expansion*
- *Backpropagation*

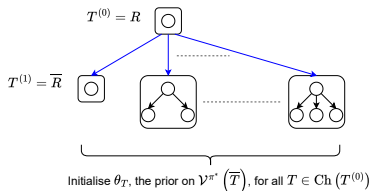
# TSDT

## The Algorithm

Initialise  $\pi(T) = \bar{T}$

**Iteration**  $t = 1$

- Selection
- Simulation
- Expansion
- Backpropagation



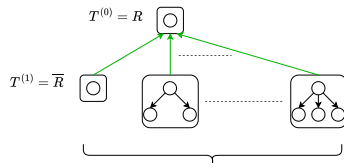
# TSDT

## The Algorithm

Initialise  $\pi(T) = \bar{T}$

**Iteration**  $t = 1$

- *Selection*
- *Simulation*
- *Expansion*
- *Backpropagation*



$$\text{Update } \theta_{T^{(0)}} \approx \max_{T \in \text{Ch}(T^{(0)})} \{-\lambda \mathbb{1}\{T \neq \bar{T}^{(0)}\} + \theta_T\}$$

$$\text{Update } \pi(T|T^{(0)}) = \mathbb{P} \left[ -\lambda \mathbb{1}\{T \neq \bar{T}^{(0)}\} + \theta_T = \max_{T' \in \text{Ch}(T^{(0)})} \{-\lambda \mathbb{1}\{T' \neq \bar{T}^{(0)}\} + \theta_{T'}\} \right]$$

### Iteration $t$

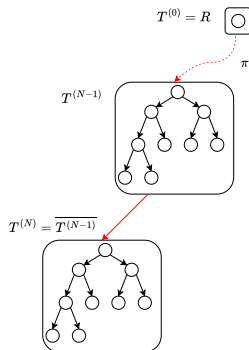
- *Selection*
- *Simulation*
- *Expansion*
- *Backpropagation*

# TSDT

## The Algorithm

### Iteration $t$

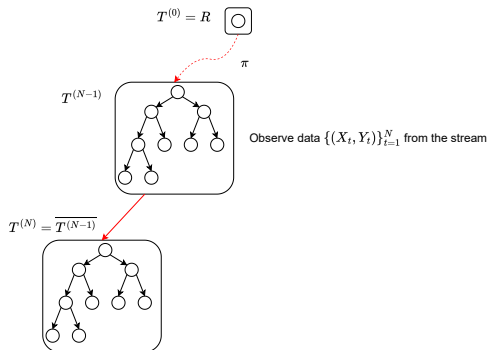
- Selection
- Simulation
- Expansion
- Backpropagation





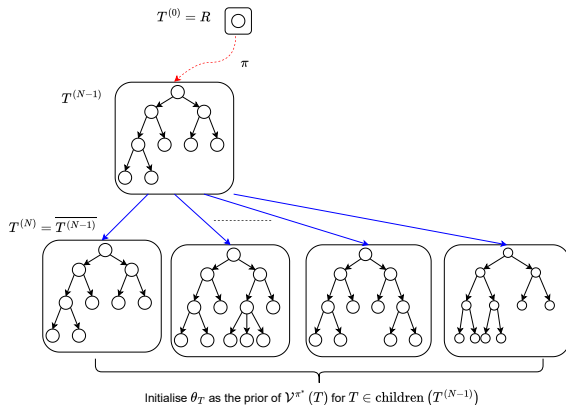
### Iteration $t$

- Selection
- Simulation
- Expansion
- Backpropagation



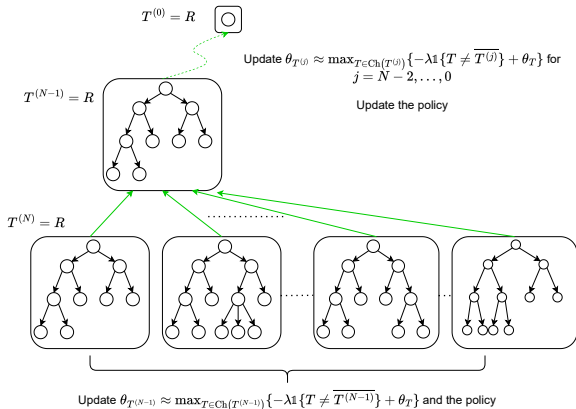
### Iteration $t$

- Selection
- Simulation
- Expansion
- Backpropagation



### Iteration $t$

- Selection
- Simulation
- Expansion
- Backpropagation



### Theorem

For any state  $T$ , TSDT satisfies the following:

$$\mu_T \xrightarrow[t \rightarrow \infty]{\text{a.s.}} \mathcal{V}^{\pi^*}(T), (\sigma_T)^2 \xrightarrow[t \rightarrow \infty]{\text{a.s.}} 0$$

For any non-terminal state  $T$ :

$$\pi_t \left( \pi^*(T) \mid T \right) \xrightarrow[t \rightarrow \infty]{\text{a.s.}} 1$$

## Section 3

# Experiments

# Experiments

## Test accuracies

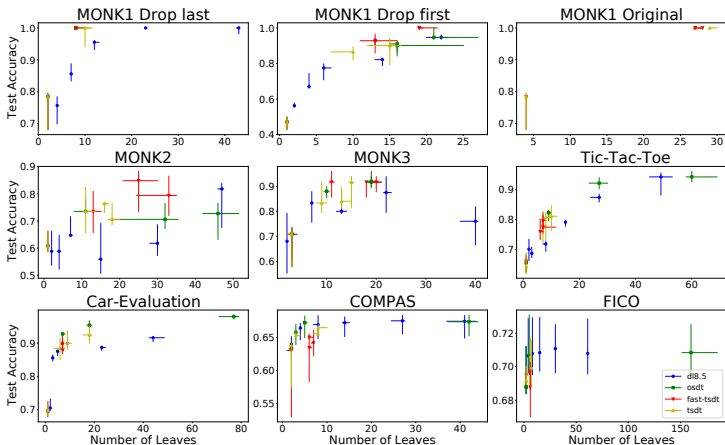


Figure: Cross-validation test accuracy comparison between TSDT, Fast-TSDT, OSDT and DL8.5 as a function of the number of leaves.

# Experiments

## Execution Times

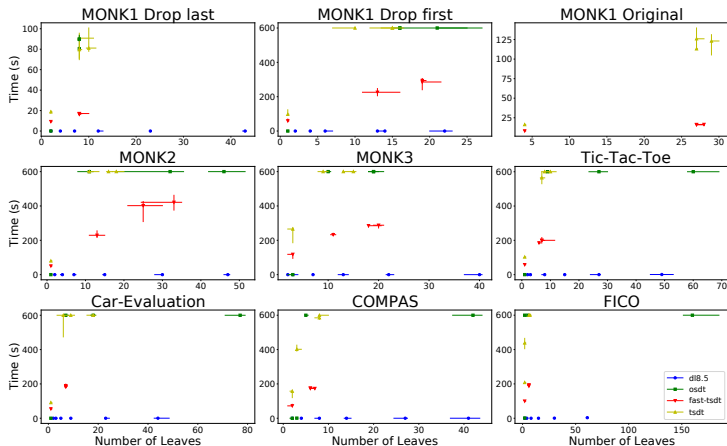


Figure: Cross-validation execution times of TSDT, OSDT and DL8.5 as a function of the number of leaves.

## Section 4

# Limitations and Future Work



- TSDT is limited to categorical features for now.
- TSDT assumes i.i.d. data, it could suffer if the data distribution is non-stationary.
- No finite-time theoretical guarantees.
- The Normal approximations of the Beta prior could cause TSDT to be suboptimal in terms of the number of iterations it takes to find  $T^*$ .

- Extend TSDT to handling numerical features.
- Extend this Value Iteration method to frequentist policies like UCB and Epsilon-greedy.
- Derive finite-time guarantees in the form of PAC-style bounds or rates of convergence.

*Thank you.*