

**UNIVERSITÉ DE LILLE**  
**INRIA**

École doctorale École Graduée MADIS-631

Unité de recherche **Centre de Recherche en Informatique, Signal et Automatique de Lille**

Thèse présentée par **Hector KOHLER**

Soutenue le **1<sup>er</sup> décembre 2025**

En vue de l'obtention du grade de docteur de l'Université de Lille et de l'Inria

**Discipline Informatique**  
**Spécialité Informatique et Applications**

# Interpretabilité via l'Apprentissage Supervisé ou par Renforcement d'Arbres de Décisions

**Thèse dirigée par** Philippe PREUX directeur  
Riad AKROUR co-directeur

## Composition du jury

<i>Rapporteurs</i>	René DESCARTES Denis DIDEROT	professeur à l'IHP directeur de recherche au CNRS	
<i>Examineurs</i>	Victor HUGO Sophie GERMAIN Joseph FOURIER Paul VERLAINE	professeur à l'ENS Lyon MCF à l'Université de Paris 13 chargé de recherche à l'INRIA chargé de recherche HDR au CNRS	président du jury
<i>Invité</i>	George SAND		
<i>Directeurs de thèse</i>	Philippe PREUX Riad AKROUR	professeur à l'Université de Lille Inria	

## COLOPHON

Mémoire de thèse intitulé « Interprétabilité via l'Apprentissage Supervisé ou par Renforcement d'Arbres de Décisions », écrit par Hector KOHLER, achevé le 30 juillet 2025, composé au moyen du système de préparation de document  $\text{\LaTeX}$  et de la classe yathesis dédiée aux thèses préparées en France.

UNIVERSITÉ DE LILLE  
INRIA

École doctorale École Graduée MADIS-631

Unité de recherche Centre de Recherche en Informatique, Signal et Automatique de Lille

Thèse présentée par **Hector KOHLER**

Soutenue le 1<sup>er</sup> décembre 2025

En vue de l'obtention du grade de docteur de l'Université de Lille et de l'Inria

Discipline **Informatique**  
Spécialité **Informatique et Applications**

# Interpretabilité via l'Apprentissage Supervisé ou par Renforcement d'Arbres de Décisions

Thèse dirigée par Philippe PREUX directeur  
Riad AKROUR co-directeur

## Composition du jury

<i>Rapporteurs</i>	René DESCARTES Denis DIDEROT	professeur à l'IHP directeur de recherche au CNRS	
<i>Examineurs</i>	Victor HUGO Sophie GERMAIN Joseph FOURIER Paul VERLAINE	professeur à l'ENS Lyon MCF à l'Université de Paris 13 chargé de recherche à l'INRIA chargé de recherche HDR au CNRS	président du jury
<i>Invité</i>	George SAND		
<i>Directeurs de thèse</i>	Philippe PREUX Riad AKROUR	professeur à l'Université de Lille Inria	



UNIVERSITÉ DE LILLE  
INRIA

Doctoral School École Gradué MADIS-631

University Department Centre de Recherche en Informatique, Signal et Automatique de Lille

Thesis defended by **Hector KOHLER**

Defended on **December 1, 2025**

In order to become Doctor from Université de Lille and from Inria

Academic Field **Computer Science**

Speciality **Computer Science and Applications**

# Interpretability through Supervised or Reinforcement Learning of Decision Trees

**Thesis supervised by** Philippe PREUX    Supervisor  
Riad AKROUR    Co-Supervisor

## Committee members

<i>Referees</i>	René DESCARTES	Professor at IHP	
	Denis DIDEROT	Senior Researcher at CNRS	
<i>Examiners</i>	Victor HUGO	Professor at ENS Lyon	Committee President
	Sophie GERMAIN	Associate Professor at Université de Paris 13	
	Joseph FOURIER	Junior Researcher at INRIA	
	Paul VERLAINE	HDR Junior Researcher at CNRS	
<i>Guest</i>	George SAND		
<i>Supervisors</i>	Philippe PREUX	Professor at Université de Lille	
	Riad AKROUR	Inria	



---

**INTERPRETABILITÉ VIA L'APPRENTISSAGE SUPERVISÉ OU PAR RENFORCEMENT D'ARBRES DE DÉCISIONS****Résumé**

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Ut purus elit, vestibulum ut, placerat ac, adipiscing vitae, felis. Curabitur dictum gravida mauris. Nam arcu libero, nonummy eget, consectetur id, vulputate a, magna. Donec vehicula augue eu neque. Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Mauris ut leo. Cras viverra metus rhoncus sem. Nulla et lectus vestibulum urna fringilla ultrices. Phasellus eu tellus sit amet tortor gravida placerat. Integer sapien est, iaculis in, pretium quis, viverra ac, nunc. Praesent eget sem vel leo ultrices bibendum. Aenean faucibus. Morbi dolor nulla, malesuada eu, pulvinar at, mollis ac, nulla. Curabitur auctor semper nulla. Donec varius orci eget risus. Duis nibh mi, congue eu, accumsan eleifend, sagittis quis, diam. Duis eget orci sit amet orci dignissim rutrum.

Nam dui ligula, fringilla a, euismod sodales, sollicitudin vel, wisi. Morbi auctor lorem non justo. Nam lacus libero, pretium at, lobortis vitae, ultricies et, tellus. Donec aliquet, tortor sed accumsan bibendum, erat ligula aliquet magna, vitae ornare odio metus a mi. Morbi ac orci et nisl hendrerit mollis. Suspendisse ut massa. Cras nec ante. Pellentesque a nulla. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus. Aliquam tincidunt urna. Nulla ullamcorper vestibulum turpis. Pellentesque cursus luctus mauris.

**Mots clés :** apprentissage par renforcement, arbres de décision, interprétabilité, méthodologie

---

**INTERPRETABILITY THROUGH SUPERVISED OR REINFORCEMENT LEARNING OF DECISION TREES****Abstract**

In this Ph.D. thesis, we study algorithms to learn decision trees for classification and sequential decision making. Decision trees are interpretable because humans can read through the decision tree computations from the root to the leaves. This makes decision trees the go-to model when human verification is required like in medicine applications. However, decision trees are non-differentiable making them hard to optimize unlike neural networks that can be trained efficiently with gradient descent. Existing interpretable reinforcement learning approaches usually learn soft trees (non-interpretable as is) or are ad-hoc (train a neural network then fit a tree to it) potentially missing better solutions.

In the first part of this manuscript, we aim to directly learn decision trees for a Markov decision process with reinforcement learning. In practice we show that this amounts to solving a partially observable Markov decision process. Most existing RL algorithms are not suited for POMDPs. This parallel between decision tree learning with RL and POMDPs solving help us understand why in practice it is often easier to obtain a non-interpretable expert policy—a neural network—and then distillate it into a tree rather than learning the decision tree from scratch.

The second contribution from this work arose from the observation that looking for a decision tree classifier (or regressor) can be seen as sequentially adding nodes to a tree to maximize the accuracy of predictions. We thus formulate decision tree induction as solving a Markov decision problem and propose a new state-of-the-art algorithm that can be trained with supervised example data and generalizes well to unseen data.

Work from the previous parts rely on the hypothesis that decision trees are indeed an interpretable model that humans can use in sensitive applications. But is it really the case? In the last part of this thesis, we attempt to answer some more general questions about interpretability: can we measure interpretability without humans? And are decision trees really more interpretable than neural networks?

**Keywords:** reinforcement learning, decision trees, interpretability, methodology

---

**Centre de Recherche en Informatique, Signal et Automatique de Lille**

Université de Lille-Campus scientifique – Bâtiment ESPRIT – Avenue Henri Poincaré – 59655 Villeneuve d'Ascq





# Sommaire

Résumé	vii
Sommaire	ix
Preliminary Concepts	1
<b>I A Difficult Problem : Direct Interpretable Reinforcement Learning</b>	<b>17</b>
1 A Formal Framework for the Reinforcement Learning of Decision Tree Policies	19
2 Limits of Direct (Deep) Reinforcement Learning of Decision Tree Policies	29
3 Understanding the Limits of Direct Reinforcement Learning	41
4 When transitions are uniform POIBMDPs are fully observable	57
<b>II An easier problem : Learning Decision Trees for MDPs that are Classification tasks</b>	<b>61</b>
5 DPDT-intro	63
6 DPDT-paper	67
7 Conclusion	89
<b>III Beyond Decision Trees : what can be done with other Interpretable Policies?</b>	<b>91</b>
8 Imitation and Evaluation	93
9 Evaluation	95

---

<b>10 Conclusion Imitation</b>	<b>111</b>
<b>Conclusion générale</b>	<b>113</b>
<b>Bibliographie</b>	<b>115</b>
<b>A Programmes informatiques</b>	<b>125</b>
<b>B Appendix I</b>	<b>127</b>
<b>Table des matières</b>	<b>135</b>

**Première partie**

**A Difficult Problem : Direct  
Interpretable Reinforcement  
Learning**



# A Formal Framework for the Reinforcement Learning of Decision Tree Policies

## 1.1 Learning Decision Tree policies

Decision tree policies offer transparency over neural network policies (cite). One can attribute an importance measure to each feature of the state for the decision of a tree policy which is harder to do for neural networks (cite).

Recently, algorithms have been developed to return decision tree policies for an MDP. Those algorithms, like any interpretable machine learning method, are either direct or indirect (cite). We propose an additional distinction amongst the direct methods : algorithms learning parametric trees and algorithms learning non-parametric trees.

Parametric trees are not “grown” from the root by iteratively adding internal or leaf nodes depending on the interpretability-performance trade-off to optimize, but are rather “optimized” : the depth, internal nodes arrangement, and state-features to consider in each nodes are fixed *a priori* and only the tested thresholds of each nodes are optimized similarly to how the weights of a neural network are optimized. As the reader might have guessed, those parametric trees are advantageous in that they can be learned with gradient descent and in the context of decision tree policies, with the policy gradient (cite). The downside of those approaches is that a user cannot know *a priori* what a “good” tree policy structure should be for a particular MDP : either the specified structure is too deep and pruning will be required after training or the tree

structure is not expressive enough to encode a good policy. Similar approaches exist in supervised learning exist where a parametric tree is fitted with gradient descent (cite). However their benefit over non-parametric trees have not been shown. When parametric trees are learned for MDPs; extra stabilizing tricks are required during training such as adaptive batch sizes (cite).

Non-parametric trees are the standard model in supervised learning (cite) and can naturally trade-off between interpretability and performances. However, specialized approaches are required since growing a tree from the root in an RL fashion is not possible. In the next section we present, to the best of our knowledge, the only direct approach to learn non-parametric decision tree policies for MDPs; Iterative Bounding Markove Decision Processes (cite).

Other more specialized approaches deal with tree policies either for specific MDPs like maze (cite) or for very small depth (cite) or when the MDP model is known (cite).

## 1.2 Example of the limitations of indirect methods

We believe this toy problem can reveal interesting properties of interpretable machine learning methods as this particular base MDP can in fact be optimally controlled by a depth-1 decision tree.

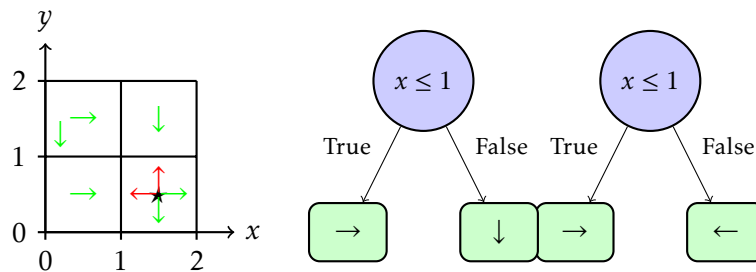


FIGURE 1.1 – Grid world MDP optimal actions. An optimal policy that takes any of the red actions in the bottom right state might not be imitated by the best depth-1 tree in terms of performance interpretability trade-offs. Center, an optimal depth-1 decision tree policy. On the right, a sub-optimal depth-1 decision tree policy.

## 1.3 Iterative Bounding Markov Decision Processes

In 2021, Topin et. al. introduced Iterative Bounding Markov Decision Processes (IBMDPs) with the promise of shifting the challenges of non-parametric decision tree

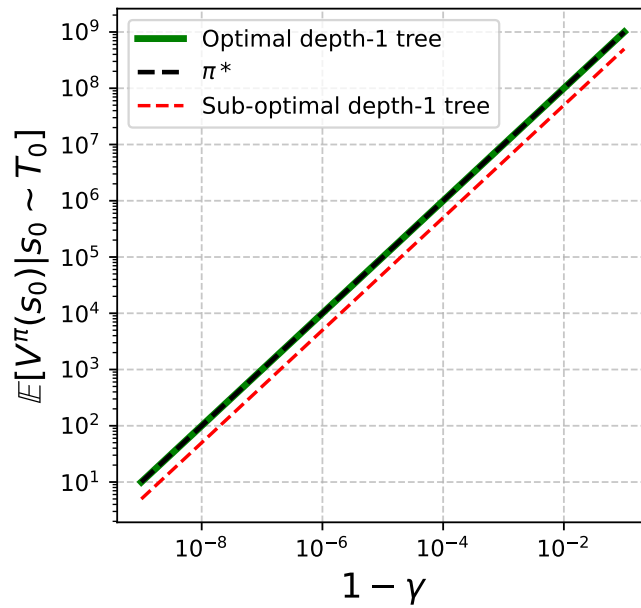


FIGURE 1.2 – MDP objective values.

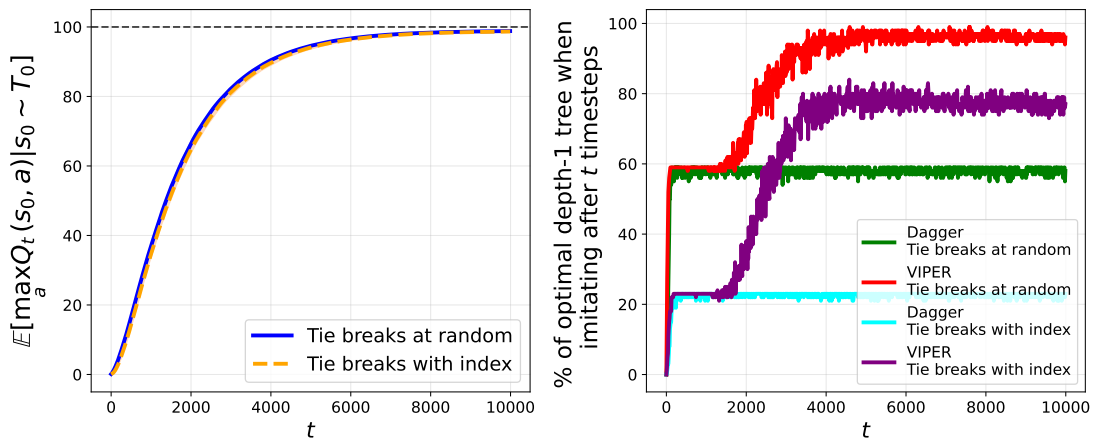


FIGURE 1.3 – Left, sample complexity curve of Q-learning with default hyperparameters on the  $2 \times 2$  grid world MDP over 100 random seeds. Right, performances of indirect interpretable methods when imitating the greedy policy with a tree at different Q-learning stages.

policy learning in the problem formulation rather than in the design of specialized reinforcement learning algorithms. Given a base MDP for which one wants to learn a decision tree policy, IBMDPs are an augmented version of this base MDP with more state features, more actions, additional reward signal, and additional transition kernel. Authors showed that certain IBMDP policies, that can be learned with RL, are equivalent to decision tree policies for the base MDP.

### 1.3.1 Formalism

The key thing to know about IBMDPs is that they are, as their name suggests, MDPs. Hence they inherit all their properties such as existence of a deterministic optimal Markovian policy. The states in an IBMDP are concatenations of base MDP states and some observations. Those observations are information about the base states that are refined—“iteratively bounded”—at each step and represent a subspace of the base MDP state space. Actions available in an IBMDP are the actions of the base MDP, that change the state of the latter, and *information gathering* actions that change the observation part of the IBMDP state. Now, taking base actions in an IBMDP is rewarded like in the base MDP, this ensures that the base objective, e.g. balancing the pole or treating cancer, is still encoded in the IBMDP reward. When taking such *information gathering* actions; the reward is an arbitrary value supposed to trade-off between performance and interpretability.

Before showing how to get decision trees from IBMDP policies, we give a formal definition of the latter following Topin et. al. (cite).

**Definition 1** (Iterative Bounding Markov decision process). *Given a factored (cite) MDP  $\mathcal{M}$  (cite), that is, a tuple  $\langle S, A, R, T, T_0 \rangle$  with  $S \subseteq \mathbb{R}^n$ , an Iterative Bounding MDP  $\mathcal{M}_{ib}$  is a tuple :*

$$\langle \overbrace{S, O}^{\text{State space}}, \underbrace{A, A_{info}}_{\text{Action space}}, \overbrace{R, \zeta}^{\text{Reward function}}, \underbrace{T_{info}, T, T_0}_{\text{Transition kernels}} \rangle$$

- $S$  the base MDP state space should be of the form  $S = [L_1, U_1] \times \dots \times [L_n, U_n]$  with  $\infty < L_i \leq U_i < \infty \forall 1 \leq i \leq n$ .
- $O$  are the observations in an IBMDP. They are partial information about the values of base MDP states :  $O \subseteq S^2 = [L_1, U_1] \times \dots \times [L_n, U_n] \times [L_1, U_1] \times \dots \times [L_n, U_n]$ . So the complete IBMDP state space is  $(S, O) = S \times O$  the concatenations of states and



observations.

- $A$  are the actions of the base MDP.
- $A_{info}$  are added information gathering actions (IGAs) of the form  $\langle i, v \rangle$  where  $i$  is a state feature index  $1 \leq i \leq n$  and  $v$  is a real number. So the complete action space of an IBMDP is the set of base actions and information gathering actions  $A \cup A_{info}$ .
- $R: S \times A \rightarrow \mathbb{R}$  is the base MDP reward function that maps base states and actions to a real-valued reward signal.
- $\zeta$  is a reward signal for taking an information gathering action. So the IBMDP reward function is to get a reward from the base MDP if the action is a base MDP action or to get  $\zeta$  if the action is a information gathering action.
- $T_{info}: S \times O \times (A_{info} \cup A) \rightarrow \Delta(S \times O)$  is the transition kernel of IBMDPs. Given the current observation  $o_t = (L'_1, U'_1, \dots, L'_n, U'_n) \in O$  and the current state is  $s_t = (s_1, s_2, \dots, s_n)$  if an AIG  $\langle i, v \rangle$  is taken, only the bounds in the observation change :

$$o_{t+1} = \begin{cases} (L'_1, U'_1, \dots, L'_i, \min\{v, U'_i\}, \dots, L'_n, U'_n) & \text{if } s_i \leq v \\ (L'_1, U'_1, \dots, \max\{v, L'_i\}, U'_i, \dots, L'_n, U'_n) & \text{if } s_i > v \end{cases}$$

If a base action  $a \in A$  is taken,  $o_{t+1}$  is reset to the default state bounds  $(L_1, U_1, \dots, L_n, U_n)$  and the base state changes according to the base MDP transition kernel:  $s_{t+1} \sim T(s, a)$ . At initialization, the base part of the IBMDP states is drawn from  $T_0$  and the observation is set always set to  $(L_1, U_1, \dots, L_n, U_n)$ . The overall IBMDP transitions are given by either  $T_{info}$ , which is fully deterministic, if an IGA is played, and by the base MDP's transition kernel otherwise.

Now remains to extract a decision tree policy for  $\mathcal{M}$  from a policy for IBMDP  $\mathcal{M}_{IB}$ .

### 1.3.2 From Policies to Trees

One can notice that *information gathering* actions resemble the Boolean functions that make up internal decision tree nodes (cite). Indeed, an agent evolving in an IBMDP essentially builds a tree by taking sequences of IGAs (internal nodes) and then a base action (leaf node) and repeats this process over time. However not all IBMDP policies are decision tree policies. In particular, only deterministic policies depending solely on the observation part of the IBMDP states  $\pi: O \rightarrow A \cup A_{info}$  are decision tree policies  $\pi_T: S \rightarrow A$  for the base MDP (cite). Algorithm (cite) from (cite) extracts a decision tree policy for the base MDP from a deterministic partially observable IBMDP policy.

**Data :** IBMDP policy  $\pi$  and observation  $o = (L'_1, U'_1, \dots, L'_n, U'_n)$   
**Result :** Decision tree policy extracted from  $\pi$   
**Function** Subtree\_From\_Policy( $o, \pi$ ) :  
      $a \leftarrow \pi(o)$   
     **if**  $a \in A_{info}$  **then**  
         **return** Leaf\_Node(action :  $a$ ) // Leaf if base action  
     **end**  
     **else**  
          $\langle i, v \rangle \leftarrow a$  // Splitting action is feature and value  
          $o_L \leftarrow o; \quad o_R \leftarrow o$   
          $o_L \leftarrow (L'_1, U'_1, \dots, L'_i, v, \dots, L'_n, U'_n); \quad o_R \leftarrow (L'_1, U'_1, \dots, v, U'_i, \dots, L'_n, U'_n)$   
          $child_L \leftarrow \text{Subtree\_From\_Policy}(o_L, \pi)$   
          $child_R \leftarrow \text{Subtree\_From\_Policy}(o_R, \pi)$   
         **return** Internal\_Node(feature :  $i$ , value :  $v$ , children : ( $child_L, child_R$ ))  
     **end**

**Algorithm 1 :** Extract a Decision Tree Policy from an IBMDP policy  $\pi$ , beginning traversal from  $obs$ .

### 1.3.3 Didactic example

We now present a didactic example of how some policies in IBMDPs correspond to a decision tree policy. Suppose a factored MDP representing a grid world with 4 cells (cite). The state space is  $S = \{(0.5, 0.5), (0.5, 1.5), (1.5, 1.5), (1.5, 0.5)\} \subseteq [0, 2] \times [0, 2]$ . The actions space are the cardinal directions  $A = \{\rightarrow, \leftarrow, \downarrow, \uparrow\}$  that shift the states by one as long as the coordinates remain in the grid. The reward for taking any action is 0 except when in the bottom right state  $s_g = (1.5, 0.5)$  which is an absorbing state : once in this state, you stay there forever.

In Figure (cite), we present how trajectories in IBMDP correspond to a decision tree policy. We build an IBMDP with two information gathering actions and nine observations.

- $\langle x, 1 \rangle$  that tests if  $x \leq 1$
- $\langle y, 1 \rangle$  that tests if  $y \leq 1$

The observations are induced by these to IGAs. The initial observation is  $o_0 = (0, 2, 0, 2)$ . The set of observations is finite and can be inferred from the possible bounds over the base states induced by the IGAs. For example when starting in state  $s_1 = (0.5, 1.5)$ , and taking first  $\langle x, 1 \rangle$  then  $\langle y, 1 \rangle$  the corresponding observations are first  $o_1 = (0, 1, 0, 2)$  and then  $o_2 = (0, 1, 1, 2)$ . The full observation set is  $O = \{(0, 2, 0, 2), (0, 1, 0, 2), (0, 2, 0, 1), (0, 1, 0, 1), (1, 2, 0, 2), (1, 2, 0, 1), (1, 2, 1, 2), (1, 2, 1, 1)\}$ . The transitions and rewards are given by definition (cite).

We summarize the goal of the next chapters in Figure (cite) : we learn decision

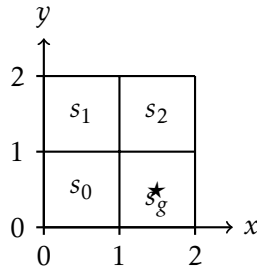


FIGURE 1.4 – A grid world MDP state space : it is a  $2 \times 2$  grid.

tree policies for some MDP  $\mathcal{M}$ , we first define an IBMDP  $\mathcal{M}_{IB}$  with some information gathering actions and their reward  $\zeta$ , then we can find the policy from observations to actions that maximizes the expected value of initial IBMDP states. We discuss in more details the impact of learn policies depending on partial observations of IBMDP in Chapter (cite) but first we attempt to reproduce the original results of (cite).

## 1.4 The promise of IBMDPs

1. Can train decision tree policy for discrete actions MDP using any recent RL algorithms.
2. Could generalize to more complex interpretable policy classes like oblique decision trees or in general any policy that can be expressed as a composition of Boolean functions.

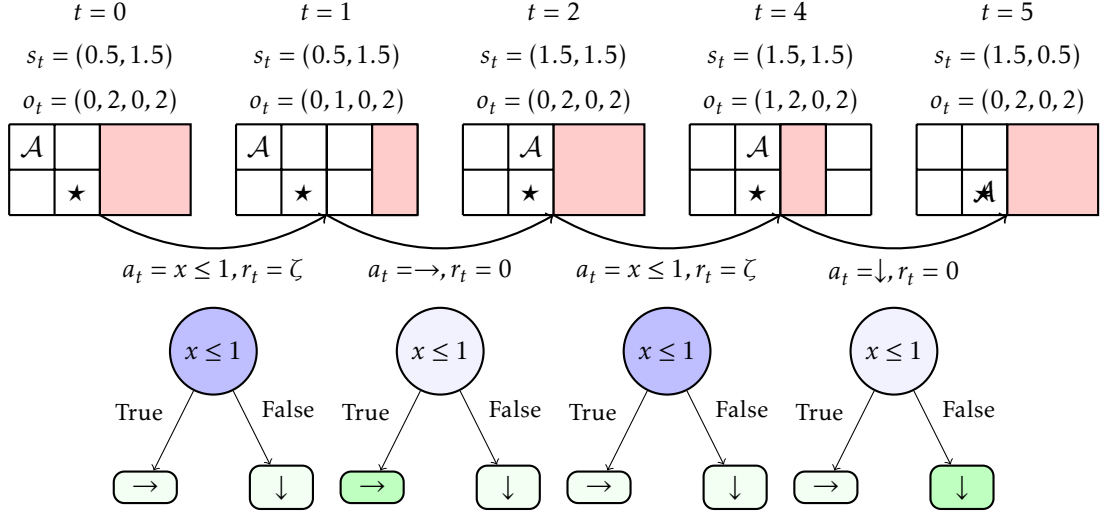


FIGURE 1.5 – An IBMDP trajectory when the base MDP is  $2 \times 2$  grid world, and equivalent decision tree policy traversal. At  $t = 0$ , the agent is in  $s_1 = (0.5, 1.5)$ . The initial observation is always the base MDP state bounds, here  $(0, 2, 0, 2)$  because the base states are in  $[0, 2] \times [0, 2]$ . The agent takes an information gathering action in  $A_{info}$  that tests the feature  $x$  of the states against the value 1. For that, the agent receives reward  $\zeta$ . This transition corresponds to going through an internal node in a decision tree policy as illustrated in the figure. At  $t = 1$ , after gathering the information that the  $x$ -value of the current base state is below 1, the observation is updated with the refined state bounds  $o_t = (0, 1, 0, 2)$  and the base state remains unchanged. The agent then takes a base action that is to move right. This gives a reward 0, reinitialized the observation to the original bounds, and changes the base state to  $s_t = (1.5, 1.5)$ . And the trajectory continues like this until the agent reaches the absorbing state  $s_g = (1.5, 0.5)$ .

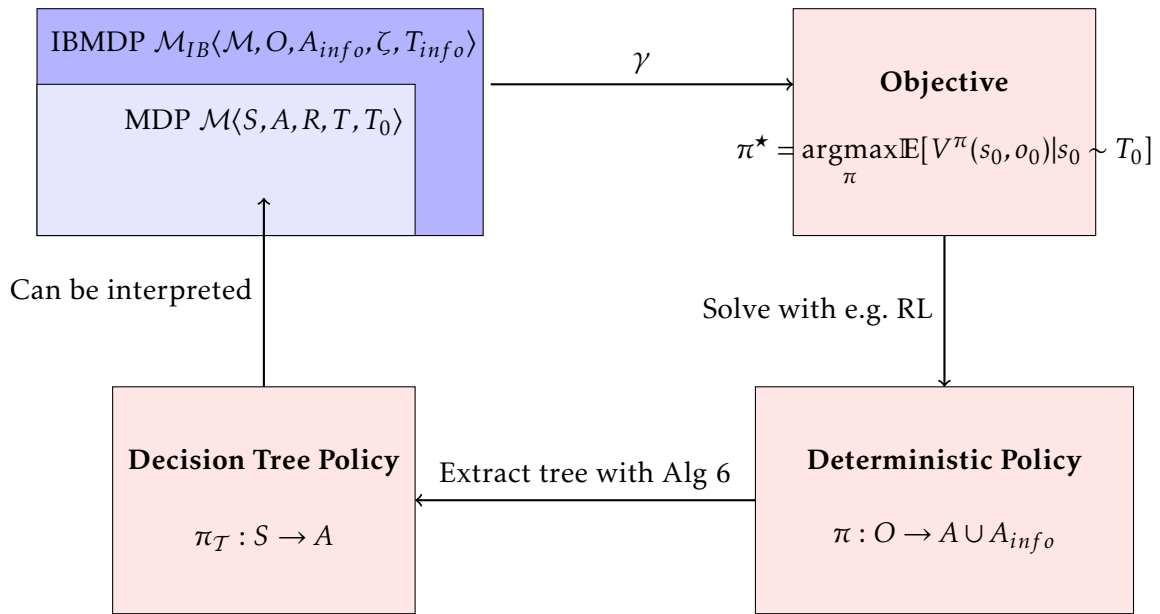


FIGURE 1.6 – A formal framework to learn decision tree policies for MDPs. This include learning a partially observable deterministic policy in a POIBMDP (cite).



# Limits of Direct (Deep) Reinforcement Learning of Decision Tree Policies

In this chapter, we reproduce the algorithms presented in (cite). In particular, we attempt to reproduce the results from Table (cite) in which authors use deep reinforcement learning to directly learn a depth-2 decision tree for the CartPole control (cite). We run the *same* algorithms using the *same* hyperparameters when possible. Surprisingly, we are unable to reproduce their results. Even after searching for different hyperparameters, we are unable to learn decision tree policies for CartPole with reinforcement learning. In this chapter we provide a first set of failure modes for the direct reinforcement learning of decision tree policies and only in the next chapter will we provide more grounded insights. (MAYBE ADD CARTPOLE FIG)

## 2.1 Reproducing “Iterative Bounding MDPs : Learning Interpretable Policies via Non-Interpretable Methods”

In (cite), authors formulate an Iterative Bounding Markov Decision Process (cite) for the learning of a depth-2 decision tree policy when the base MDP is CartPole (cite). Then they use reinforcement learning to learn a partially-observable policy from which they can extract a decision tree using Alg 6 (cite).

Authors propose two tricks to improve learning : 1) limit the tree depth, and 2), parametric information gathering actions that adapt to the state.

### 2.1.1 IBMDP formulation

As described in (cite), given a base MDP  $\mathcal{M}\langle S, A, R, T, T_0 \rangle$ , to define an IBMDP  $\mathcal{M}_{IB}\langle S, O, A, A, R, \zeta, T, T_0, T_{info} \rangle$  given a base MDP, the user needs to provide the set of information gathering actions  $A_{info}$  and the reward  $\zeta$  for taking those. The observation space of the IBMDP  $O$  and the deterministic transitions  $T_{info}$  are induced from the given hyperparameters  $A_{info}, \zeta$ .

Authors propose a to parametrize the set of IGAs with  $i \times p$  actions  $\langle v_k, i \rangle$  with  $v_k$  depending on the current observation  $o = (L'_1, U'_1, \dots, L'_i, U'_i, \dots, L'_n, U'_n) : v_k = \frac{k(U'_i - L'_i)}{p+1}$ . This parametric IGAs space keeps the discrete IBMDP action space at a reasonable size while providing a learning agent with very varied tests.

For example, if we define an IBMDP with  $p = 3$  for the grid-world MDP of the previous chapter, at every time step the agent can take one of six actions. At  $t = 0$ , recall that  $o_0 = (0, 2, 0, 2)$ , so if an agent takes, e.g., IGA  $\langle v_2, 2 \rangle$ , the effective IGA is  $\langle v_2 = \frac{k(2-0)}{3+1}, i \rangle = \langle 1, 2 \rangle$  which in turn effectively corresponds to a test node “ $y \leq 1$ ?” for a decision tree policy acting in the grid-world. If the base state has  $y$  value 0.5, then the next observation at  $t = 1$  is  $o_1 = (0, 2, 0, 1)$ . Now, if the agent were to take the same IGA  $\langle v_2, 2 \rangle$  again, it would be effectively  $\langle v_2 = \frac{k(1-0)}{3+1}, i \rangle = \langle 0.5, 2 \rangle$ . This would give the next observation at  $t = 2$   $o_2 = (0, 2, 0, 0.5)$  and so on  $\dots$

Furthermore, akin to the supervised learning setting (cite), author propose to regularize the learned decision tree policy with a maximum depth parameter  $D$ . Unfortunately, authors did not describe as they implemented the depth control in their work, hence we have to try different approaches to reproduce their results.

To control the tree depth during learning we can either penalize the agent everytime it takes more than  $D$  IGAs since the last base action, or we could send a termination signal to the agent. This approach interrupts the trajectory : the agent stops acting. However the latter, in theory, breaks the MDP formalism as a termination signal is different from an absorbing state. The penalization approaches can also break the MDP formalism because the reward function now depends on time. We will try both approaches in our experiments. We actually find that when  $p+1$ , the IBMDP information gathering space parameter, is a prime number, then as a direct consequence of the *Chinese Remainder Theorem* (cite)(proof), the current tree depth is directly encoded in the current observation  $o_t$ . Hence, when  $p+1$  is prime, we can control the depth through transitions or rewards without tracking the time.



### 2.1.2 Modified Deep Reinforcement Learning algorithms

Authors of (cite) use two deep reinforcement learning baselines to which they apply some modifications in order to learn partially-observable policies from feature bounds to information gathering actions and base actions. The first algorithm is the proximal policy optimization algorithm (PPO)(cite)(algo). This algorithm can be seen as a deep learning version of the Policy Gradient Algorithm (cite). The key modification to learn policies in IBMDPs that are equivalent to trees is to train a neural network policy  $O \rightarrow A \cup A_{info}$  rather than a policy  $S \times O \rightarrow A \cup A_{info}$  like in the traditional MDP setting. The value function is still approximated by a neural network  $S \times O \rightarrow A \cup A_{info}$ .

The second deep reinforcement learning algorithm used is the deep Q-networks algorithm (DQN) (cite); a deep learning version of the Q-learning algorithm Alg (cite). A similar modification is done to DQN to work with partially observable policy. The trained Q-function is approximated with a neural network  $O \rightarrow \mathbb{R}^{|A \cup A_{info}|}$  rather than  $S \times O \rightarrow \mathbb{R}^{|A \cup A_{info}|}$ . In this modified DQN, the temporal difference error target for the Q-function  $O \rightarrow A \cup A_{info}$  is approximated by a neural network  $S \times O \rightarrow A \cup A_{info}$  that is in turn trained by bootstrapping the temporal difference error with itself.

Those two variants of DQN and PPO have first been introduced in (PINTO) for robotics task and later studied theoretically to solve POMDPs (cite) in Baisero's work. Baisero coined those approaches asymmetric reinforcement learning for POMDPs as an agent accesses full state information during the learning of a partially observable policy, hence using "asymmetric" information. We study the connections between, learning decision tree policies for MDPs, and, POMDPs, extensively in the next chapter but for now we focus on reproducing the original results from (cite) in which POMDPs and asymmetric RL were not mentioned.

Next, we present the precise experimental setup we use to reproduce the work of (cite) in order to study deep reinforcement learning of decision tree policies for the CartPole MDP.

## 2.2 Experimental setup

### 2.2.1 (IB)MDP

We use the exact same MDP and associated IBMDP for our experiments as authors of (cite) except mentioned otherwise. Next, we describe in details the MDP and IBMDP used.

**MDP** The problem at hand is the CartPole MDP. At each time step a learning agent observes the cart position velocity and the pole angle and angular velocity, and can take action to push the cart left or right. While the cart is roughly balanced, i.e., while the cart angle remains in some fixed range, the agent gets a positive reward. If the cart is out of balance; the agent goes to an absorbing terminal state and gets 0 reward forever. In practice, authors use the gymnasium `CartPole-v0` implementation (cite) of the CartPole MDP in which trajectories are truncated after 200 timesteps making the maximum cumulative reward over a trajectory to be 200. Since the IBMDP definition requires the MDP state space to be a factor of bounded segments, authors bound the the state space of the CartPole MDP to be in  $[-2, 2] \times [-2, 2] \times [-0.14, 0.14] \times [-0.14, 0.14]$ .

**IBMDP** Authors define the associated IBMDP with  $\zeta = -0.01$  and parametric information gathering action space defined by  $p = 3$ . In addition we also try  $\zeta = 0.01$  and  $p = 2$ . The discount factor used by the authors is  $\gamma = 1$ .

We potentially differ from the original paper setting in the way we handle maximum depth limitation. Indeed authors restrain the learning of policies to be equivalent to depth-2 trees but don't detail how they do so. We hence try two different approaches as mentioned in the previous section : terminating trajectories if the agent takes too much information gathering in a row or simply giving a reward of  $-1$  to the agent everytime it takes an information gathering action past the depth limit. This makes the transition or the reward function non markovian respectively as those operations are time dependent.

We will also try IBMDPs where we do not limit the maximum depth for completeness.

## 2.2.2 Baselines

**Modified DQN** as mentioned above, authors use an asymmetric version (cite) of the DQN algorithm (cite). We use the exact same hyperparameters for DQN as the authors when possible. We use the same layers width (128) and number of hidden layers (2), the same exploration strategy ( $\epsilon$ -greedy with linearly decreasing value  $\epsilon$  between 0.5 and 0.05 during the first 10% of the training), the same replay buffer size ( $10^6$ ) and the same number of transitions to be collected randomly before doing value updates ( $10^5$ ). We also try to use more exploration during training (change the initial  $\epsilon$  value to 0.9). We use the same optimizer (RMSprop (cite) with hyperparameter 0.95 and learning rate  $2.5 \times 10^{-4}$ ) to update the Q-networks.

Authors did not share what DQN implementation they used so we used the stable-baselines3 one (cite). Authors did not share what activations they used so we try both

$\tanh()$  and  $\text{relu}()$ .

**Modified PPO** for the modified PPO algorithm, we can exactly match the authors hyperparameters since they use the open source stable-baselines3 implementation of PPO.

Similarly to authors of (cite) we train DQN on 1 million timesteps and PPO on 4 million timesteps. For PPO, authors originally used 3.2 million timesteps but we rounded above : in practice this change should only be beneficial.

**DQN and PPO** in addition to the above modified RL algorithms that learn a partially observable policy in an IBMDP equivalent to a decision tree policy for the CartPole MDP, we learn neural network policies  $S \times O \rightarrow A \cup A_{info}$  and neural network  $Q$ -functions  $S \times O \rightarrow \mathbb{R}^{|A \cup A_{info}|}$  using standard DQN and PPO on IBMDP. We also apply those standard DQN and PPO to learn neural networks policies  $S \times O \rightarrow A \cup A_{info}$  and neural network  $Q$ -functions  $S \times O \rightarrow \mathbb{R}^{|A \cup A_{info}|}$  for the base CartPole MDP. We can re-use the original paper (cite) hyperparameters to train the standard PPO and DQN.

**Indirect methods** we also learn decision tree policies by imitating the neural experts returned by standard PPO and DQN at different training stages. In particular, when fitting decision trees to the greedy policy from  $Q$ -functions returned by DQN, we can use the VIPER variant of Alg (cite). This rescales each transition  $(s, a)$  given by the expert by their importance in the future values  $|\max_{a'} Q(s, a') - \max_{a''} Q(s, a'')|$  (see section of cite). When fitting policies returned by PPO, unless we want to train an additional neural network to estimate  $Q$ -values, we are bound to use the Dagger (cite) algorithm (Alg (cite)) since PPO only returns a policy and a state-value function.

For each indirect method, we fit decision trees using the greedy tree classifier from scikit-learn with default hyperparameters and maximum depth of 2 to 10000 transitions collected by the neural experts.

We summarize hyperparameters for the IBMDP and for the learning algorithms in Tables (cite) and (cite).

For standard RL on IBMDP, we instantiate the IBMDP with original hyperparameters and depth control with negative rewards and we instantiate the RL algorithms with original hyperparameters and  $\tanh()$  activations. Those two unique combinations are 20 times each. For standard RL on the base MDP, we instantiate the standard RL with original hyperparameters and  $\tanh()$  to.

TABLEAU 2.1 – IBMDP hyperparameters. We try 12 different IBMDPs. In **green** we highlight the hyperparameters from the original paper and in **red** we highlight the hyperparameter names for which author do not give information.

Hyperparameter	Values
Discount factor $\gamma$	<b>1</b>
Information gathering actions parameter $p$	<b>2, 3</b>
Information gathering actions rewards $\zeta$	<b>-0.01, 0.01</b>
<b>Depth control</b>	Done signal, negative reward, none

TABLEAU 2.2 – (Modified) DQN trained on  $10^6$  timesteps. This gives four different instantiation of (modified) DQN. Hyperparameters not mentioned are stable-baselines3 default. In **green** we highlight the hyperparameters from the original paper and in **red** we highlight the hyperparameter names for which author do not give information.

Hyperparameter	Values
Buffer size	<b><math>10^6</math></b>
Random transitions before learning	<b><math>10^5</math></b>
Epsilon start	<b>0.9, 0.5</b>
Epsilon end	<b>0.05</b>
Exploration fraction	<b>0.1</b>
Optimizer	<b>RMSprop (<math>\alpha = 0.95</math>)</b>
Learning rate	<b><math>2.5 \times 10^{-4}</math></b>
Networks architectures	<b>[128, 128]</b>
<b>Networks activation</b>	tanh(), relu()

### 2.2.3 Metrics

The key metric of this section is performance when controlling the CartPole after  $t$  steps of learning. Depending on what the underlying task is, we obtain this value during the training differently for each type of learning agent. For modified RL algorithms that learn a partially observable policy (or Q-function) in an IBMPD, we periodically extract the policy (or Q-function) and use Alg 6 (cite) to extract a decision tree for the CartPole MDP. We then evaluate the tree on 100 independent trajectories in the MDP and report the mean cumulative reward. For RL applied to IBMDPs, since we can't deploy learned policies directly to the base MDP as the state dimensions mismatch, we evaluate those policy periodically to control an other IBMDP for the CartPole MDP where  $\zeta = 0$  ensuring that the cumulative IBMDP reward only corresponds to the reward for taking base actions. Similarly, we do 100 trajectories of the extracted policies in the

TABLEAU 2.3 – (Modified) PPO trained on  $4 \times 10^6$  timesteps. This gives two different instantiation of (modified) PPO. Hyperparameters not mentioned are stable-baselines3 default. In green we highlight the hyperparameters from the original paper and in red we highlight the hyperparameter names for which author do not give information.

Hyperparameter	Values
Steps between each policy gradient steps	512
Number of minibatch for policy gradient updates	4
Networks architectures	[64, 64]
Networks activations	tanh(), relu()

evaluation IBMDP and report the average cumulative reward.

For RL applied directly to the base MDP we can just periodically extract the learned policies and evaluate them on 100 trajectories of the base MDP.

Since imitation learning baselines do not train as functions of time steps, we fit trees to neural experts trained for different duration on the base MDP. We then evaluate the imitated decision tree policy on ten trajectories on the CartPole base MDP.

Every single combination of IBMDP and modified RL instantiation is run 20 times. Next, we present our results and discuss the reproducibility and limitations of the original approach presented in (cite).

## 2.3 Results

### 2.3.1 Can Modified Deep Reinforcement Learning consistently learn good decision tree policies for CartPole?

On Figure (cite) we observe that learning curves for Modified DQN on IBMDPs can learn—the curves have an increasing trend—but find poor decision tree policies for CartPole in average—the curves flatten at the end of the x-axis and have low y-values—. In, particular, among all the learning curves that could correspond to the original paper’s proposed instantiations of modified RL baselines (cite), the best learning curves in average converges to very poor CartPole tree policies. On Figure (cite) we observe that Modified PPO learn better decision tree policies than Modified DQN but this could be because we trained longer like in the original paper. The best Modified PPO that could correspond to the original paper’s instantiation is again not the best performer.

On the other hand, on Figures (cite) and (cite), we observe that baselines seeking fully observable policies (RL + IBMDP and RL + MDP), learn better CartPole policies in

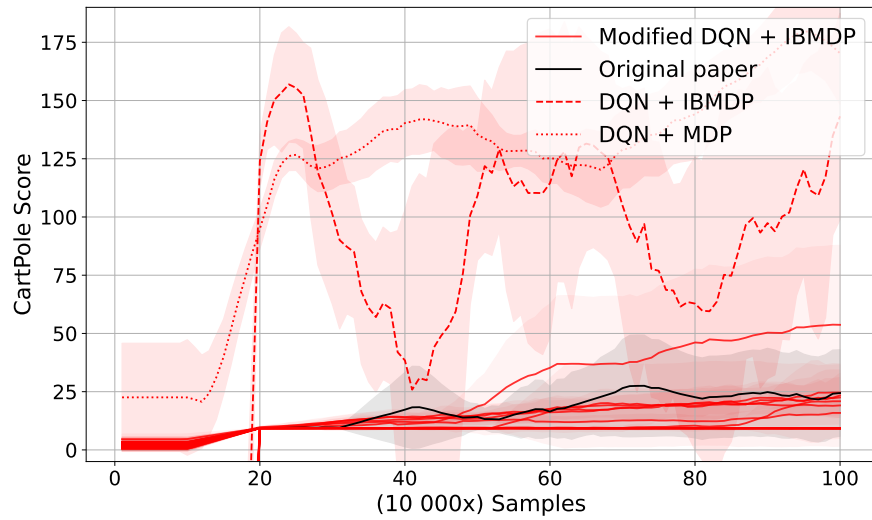


FIGURE 2.1 – Variations of Modified DQN, DQN, on CartPole IBMDP variations. We give different line-styles for the learning curves for DQN applied directly on CartPole and DQN applied on the IBMDP. Since there are multiple possible candidates for the original paper hyperparameters, we choose to color the (Modified DQN variant, IBMDP variant) pair that resulted in the best decision tree policy on CartPole among the instances that could match the original paper. Shaded areas represent the confidence interval at 95% at each measure on the y-axis.

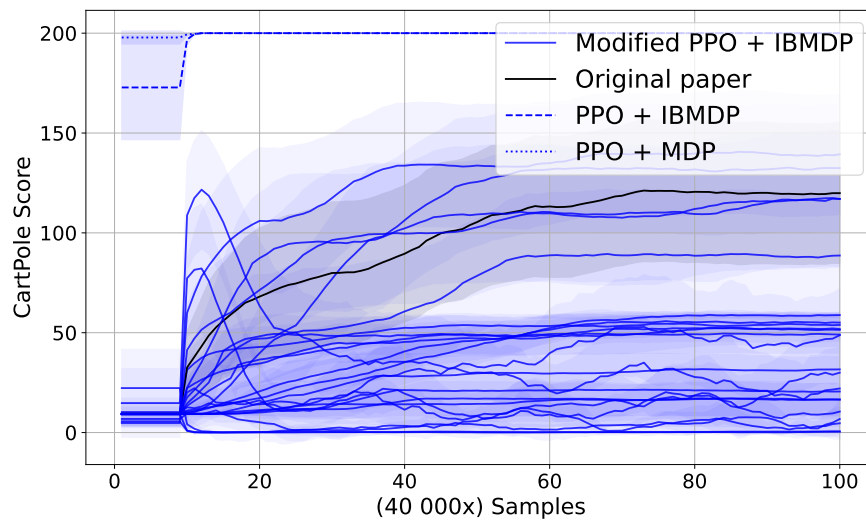


FIGURE 2.2 – Variations of Modified PPO, PPO, on CartPole IBMDP variations. We give different line-styles to the learning curves for PPO applied directly on CartPole and PPO applied on the IBMDP. Since there are multiple possible candidates for the original paper hyperparameters, we choose to color the (Modified PPO variant, IBMDP variant) pair that resulted in the best decision tree policy on CartPole among the instances that could match the original paper. Shaded areas represent the confidence interval at 95% measure on the y-axis.

average for both DQN and PPO-like baselines. We do notice that for DQN-like baselines, learning seems difficult in general independently of the setting while on Figure (cite) it is clear that for PPO baselines seeking fully observable, learning is super efficient and agent find optimal policies with reward 200.

TABLEAU 2.4 – Top 5 Hyperparameter Configurations for Modified DQN + IBMDP

Rank	$p$	Depth control	Activation	Exploration	$\zeta$	Final Performance
1	3	termination	tanh()	0.9	0.01	53
2	2	termination	tanh()	0.5	-0.01	24
3	3	termination	tanh()	0.5	-0.01	24
4	2	termination	tanh()	0.5	0.01	23
5	2	termination	tanh()	0.9	-0.01	22

TABLEAU 2.5 – Top 5 Hyperparameter Configurations for Modified PPO + IBMDP

Rank	$p$	Depth Control	Activation	$\zeta$	Final Performance
1	3	reward	relu()	0.01	139
2	3	done	relu()	0.01	132
3	3	reward	tanh()	-0.01	119
4	3	reward	relu()	-0.01	117
5	3	reward	tanh()	0.01	116

### 2.3.2 How does direct interpretable reinforcement learning perform compared to the indirect approach?

On Figure (cite), we isolate the best performing baselines that learn decision tree policies for CartPole using either Modified DQN or Modified PPO and compare them with imitation learning baselines that use the surrogate objective (cite) to find CartPole decision tree policies. We find that despite having poor performances in *average* (c.f. Figure (cite)), the deep reinforcement learning baselines find very good decision tree policies and also extremely poor ones as shown by the min-max shaded areas on the left of Figure (cite) and the histograms of final trees performances on Figure (cite). In fact, when compared with imitation learning baselines, that fit the policy of a PPO + MDP or fit the greedy policy w.r.t to the Q-network of a DQN + MDP, when running e.g., 20 seeds of direct reinforcement learning baselines, there are some seeds that can find good decision tree policies as fast as the indirect imitation learning baselines : if we look at the shaded area of Modified DQN on Figure (cite), we can observe that it is



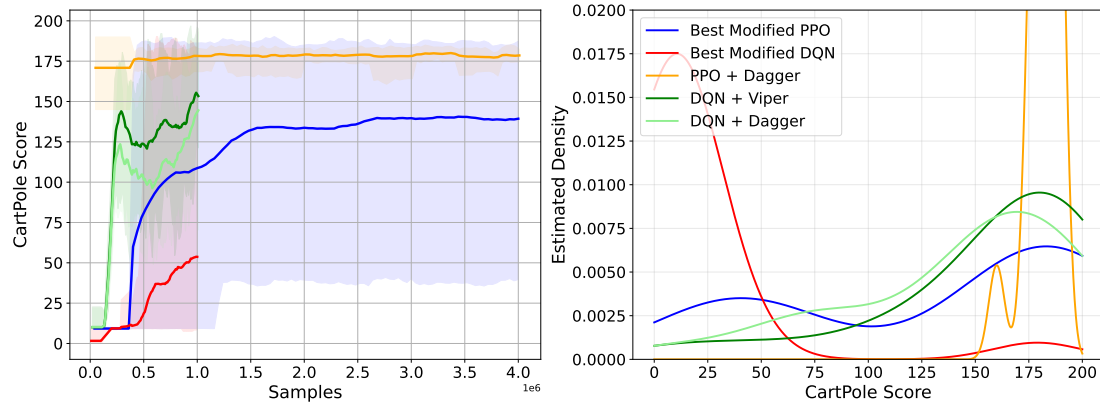


FIGURE 2.3 – (left) Mean performance of the best Modified PPO on the best IBMDP with shaded areas representing the min and max performance over the 20 seeds during training. (right) Histogram of the final decision tree policies performances.

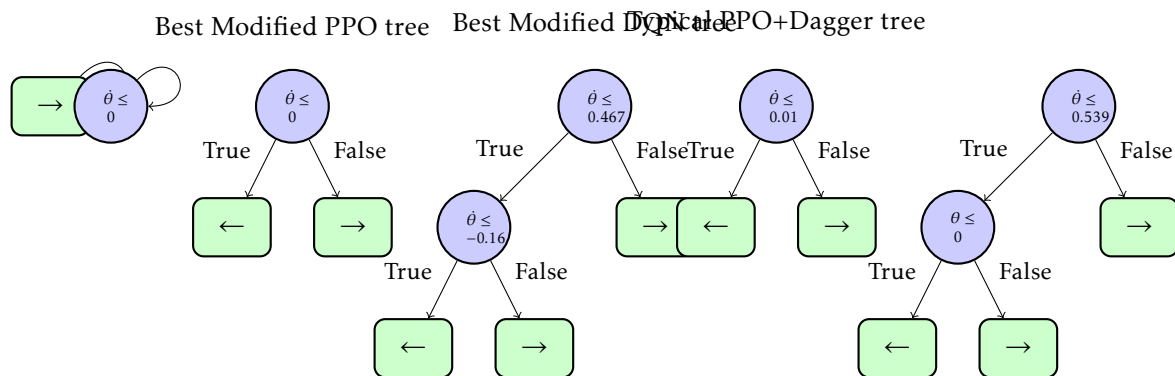


FIGURE 2.4 – Trees obtained by DRL against trees obtained with imitation.

consistently overlapping with an indirect method curve. However this is not desirable, a user want an agent that can consistently find good decision tree policy. As shown by the histogram (cite), indirect methods consistently find good decision tree policies (there is only 1 histogram bin where good CartPole scores are). On the other hand, the final trees returned by direct RL methods seem equally distributed on both extremes of the scores.

On Figure (cite), we compare the best trees returned by Modified DQN and Modified PPO. We used Algorithm 6 to extract 20 trees from the 20 partially observable policies returned by deep reinforcement learning over 20 seeds. We then plot the best tree for each baseline. Those trees get an average reward of roughly 175. Similarly, we imitation learning returns 20 trees. However unlike for direct methods, the tree returned are extremely similar across seed. In particular they often only vary in the scalar value used in the root node but in general have the same structure and test the angular velocity. On the other hand the most frequent tree across seed returned by modified RL baselines is the “infinite” tree that corresponds to only taking an information gathering action.

## 2.4 Discussion

We have shown that compared to learning non-interpretable policies for both the base MDP and the associated IBMDP, learning partially observable policies in IBMDP is very difficult even when learning policies corresponding to depth-2 decision trees for the simple CartPole task. In the next chapter we down scale the learning algorithms to highlight the connexions between the above challenges and POMDP hardness results.

# Understanding the Limits of Direct Reinforcement Learning

## 3.0.1 Partially Observable IBMDPs

Now we know that to find a decision tree policy for a given MPD  $\mathcal{M}$  satisfying the definition (cite); one has to find a partially observable—sometimes called *memoryless* or *reactive*—deterministic policy for an IBMDP  $\mathcal{M}_{IB}$ . Such problems are classical instances of Partially Observable Markov Decision Processes (POMDPs) (cite). This connection with POMDP was not done by the authors of IBMDPs.

**Definition 2** (Partially Observable Markov Decision Processes). *A Partially Observable Markov Decision Process (POMDP) is a tuple  $\langle X, A, O, T, T_0, \Omega, R \rangle$  where :*

- $X$  is the state space (like in the base definition of MDPs).
- $A$  is a finite set of actions (like in the base definition of MDPs).
- $O$  is a set of observations.
- $T : X \times A \rightarrow \Delta X$  is the transition kernel, where  $T(s, a, x') = P(x'|x, a)$  is the probability of transitioning to state  $x'$  when taking action  $a$  in state  $x$
- $T_0$  : is the initial distribution over states.
- $\Omega : X \rightarrow \Delta O$  is the observation kernel, where  $\Omega(x', a, o) = P(o|x', a)$  is the probability of observing  $o$  in state  $x$
- $R : X \times A \rightarrow \mathbb{R}$  is the reward function, where  $R(x, a)$  is the immediate reward for taking action  $a$  in state  $x$

Note that  $\langle X, A, R, T, T_0 \rangle$  defines an MDP (cite).

We can simply extend the definition of Iterative Bounding MDPs (cite) with an

observation kernel to get Partially Observable IBMDPs :

**Definition 3** (Partially Observable Iterative Bounding Markov Decision Processes). *a Partially Observable Iterative Bounding Markov Decision Process (POIBMDP) is a an IBMDP (cite) extended with an observation kernel*

$$\langle \underbrace{S, \underbrace{O}_{\text{Observations}}}_{\text{fully observable states } X}, \underbrace{A, A_{info}}_A, \underbrace{R, \zeta}_R, \underbrace{T_{info}, T, T_0}_{T, T_0}, \Omega \rangle$$

The sole specificity of POIBMDPs compared to the general definition of POMDPs, is that  $\Omega(o|(s, o'))$ , the probability of observing  $o$  in  $(s, o')$ , is  $1_{o=o'}$ . This particular instance of POMDPs with observations being some indices of the fully-observable states has other names in the litterature : Mixed Observability MDPs (cite), Block MDPs (cite). POIBMDPs can also be seen as non-stationary MDPs in which there is one different transition kernel per base MDP state : these are called Hidden-Mode MDPs (cite).

Following (cite) we can write the definition of the value of a deterministic partially observable policy  $\pi : O \rightarrow A$  in observation  $o$ .

**Definition 4** (Partial observable value function). *In a POIBMDP, the expected cumulative discounted reward of a deterministic partially observable policy  $\pi : O \rightarrow A \cup A_{info}$  starting from observation  $o$  is  $V^\pi(o)$  :*

$$V^\pi(o) = \sum_{(s, o') \in S \times O} P^\pi((s, o')|o) V^\pi((s, o'))$$

with  $P^\pi((s, o')|o)$  the asymptotic occupancy distribution (see cite for definition) of the fully observable state  $(s, o')$  given the partial observation  $o$  and  $V^\pi((s, o'))$  the classical state-value function defined in (cite).

The asymptotic occupancy distribution is the probability of a policy  $\pi$  to be in  $(s, o')$  while observing  $o$  and having taken actions given by  $\pi$ .

The problem that we solve is to find the deterministic partially observable policy that maximizes the excpeted value in the initial observation :

$$\pi^\star = \operatorname{argmax}_{\pi} J(\pi) = \operatorname{argmax}_{\pi} V^\pi(o_0) \quad (3.1)$$

With  $\pi : O \rightarrow A \cup A_{info}$ . There is no expectation over possible initial observation in the above objective function as there is only one initial observation in a POIBMDP :

$$o_0 = (L_1, U_1, \dots, L_n, U_n).$$

This particular problem of learning deterministic partially observable policies for POMDPs has been studied in the works of Littman, Singh and Jordan : (cite). In (cite) atuhors give som intuition behind why the above optimization problem is hard. For example, the optimal partially observable policy can be stochastic (cite precise section), hence policy gradient algoiriothms (cite) are to avoid. Furthermore, the optimal deterministic patially observable policy might not maximize all the value of all observations simulataneously (cite precise section) which makes difficult the use the Bellman optimality equation (cite) to compute policies. Despite those hardness results, empirical results of applying RL to POMDPs by naively setting the states to be observations has shown promising results (cite). More recently, the framework of Baisero et. al. (cite) called asymmetric RL has also shown promising empirical and theoretical results when leveraging fully-observable state information during training of a partially observable policy. In the next chapter, we use reinforcement learning to train decision tree policies for MDPs by seeking deterministic partially observable policies for POIBMDPs (cite).

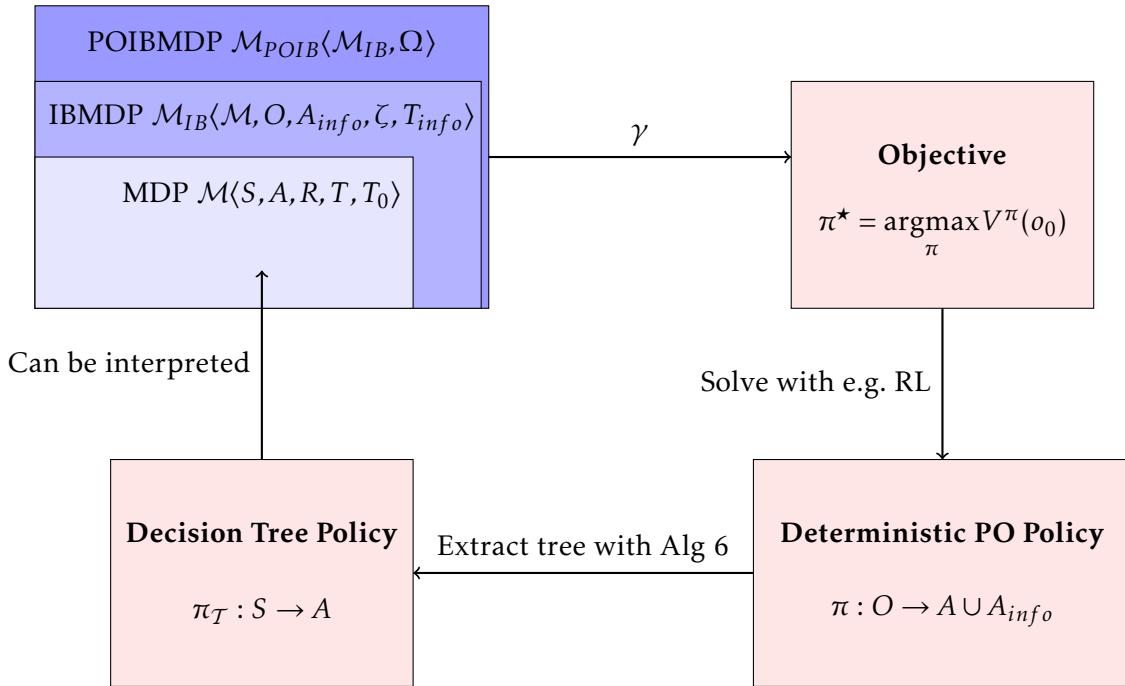


FIGURE 3.1 – A formal framework to learn decision tree policies for MDPs. This include learning a partially observable deterministic policy in a POIBMDP (cite).

In this chapter we investigate the properties of reinforcement learning for learning a depth-1 decision tree policy for a very simple Markov decision process. Using the formal

framework of partially observable iterative bounding Markov decision processes (cite), we aim to learn partially observable deterministic policies that correspond to decision tree policies (cite).

We will attempt to learn a depth-1 tree policy for the  $2 \times 2$  grid world from (cite) as the base MDP by solving the example POIBMDP from the previous chapter. To focus on reinforcement learning properties in POIBMDPs rather than, e.g., modelling properties, we will carefully craft POIBMDPs such that the *optimal* partially observable deterministic policy corresponds to a depth-1 tree policy. To do so, we present next how the POIBMDP objective values (cite) from different decision tree policies change with  $\zeta$  for some fixed discount factor  $\gamma$ .

### 3.1 Constructing POIBMDPs which optimal solutions are the depth-1 tree

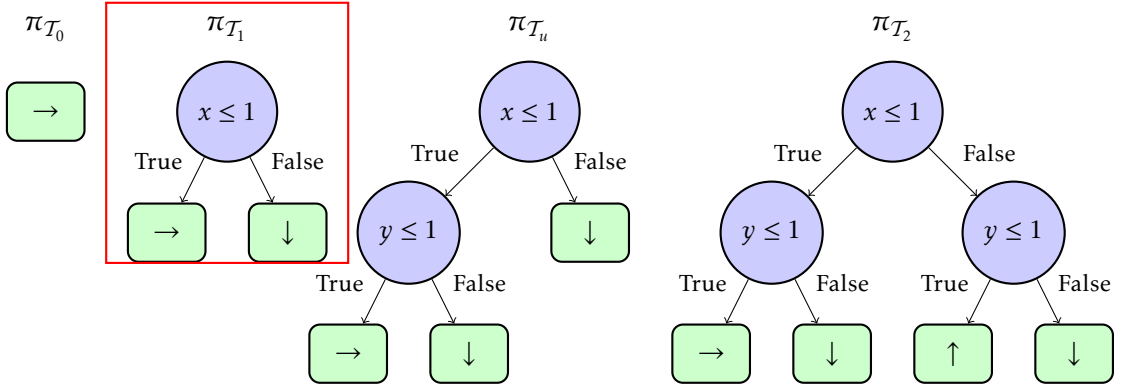


FIGURE 3.2 – For each decision tree structure, e.g., depth-1 or unbalanced depth-2, in the space of deterministic partially observable POIBMDP policies, we illustrate the decision tree which gets the highest base rewards.

In this section we compute the *exact* objective values of POIBMDP policies as functions of  $\zeta$ . After doing so, we can look for  $\zeta$  values such that the optimal POIBMDP policy corresponds to a depth-1 decision tree for the grid world MDP from the previous chapter. Indeed, because we know all the base states, all the observations, all the actions, all the rewards and all the transitions, we can compute exactly the values of different partially observable deterministic policies given  $\zeta$  the reward for IGAs and  $\gamma$  the discount factor.

Each of those policies can be one of the following trees illustrated in Figure (cite) :

- $\pi_{\mathcal{T}_0}$  : a depth-0 tree equivalent to always taking the same base action
- $\pi_{\mathcal{T}_1}$  : a depth-1 tree equivalent alternating between an IGA and a base action
- $\pi_{\mathcal{T}_u}$  : an unbalanced depth-2 tree that sometimes takes two IGAs then a base action and sometimes a an IGA then a base action
- $\pi_{\mathcal{T}_2}$  : a depth-2 tree that alternates between taking two IGAs and a base action
- an infinite “tree” that only takes IGAs

Furthermore, because from (cite) we know that for POMDPs, stochastic policies can sometimes be optimal, we also compute the value of the stochastic policy that alternates between two base actions :  $\rightarrow$  and  $\downarrow$ . Those two base actions always lead to the goal state in expectation.

We detail the calculations for the detpth-1 decision tree objective value (cite) and defer the calculations for the other policies to the Appendix (cite).

**Proposition 1** (Depth-1 decision tree objective value). *The objective value of the best depth-1 decision tree from Figure (cite) is  $V^{\pi_{\mathcal{T}_1}}(o_0) = \frac{4\zeta + \gamma + 2\gamma^3 + \gamma^5}{4(1-\gamma^2)}$ .*

*Démonstration.*  $\pi_{\mathcal{T}_1}$  has one root node that tests  $x \leq 1$  (respectively  $y \leq 1$ ) and two leaf nodes  $\rightarrow$  and  $\downarrow$ . To compute  $V^{\pi_{\mathcal{T}_1}}(o_0)$ , we compute the values of  $\pi_{\mathcal{T}_1}$  in each of the possible startin states  $(s_0, o_0), (s_1, o_0), (s_2, o_0), (s_g, o_0)$  and compute the expectation over those. At inititalization, when the base state is  $s_g = (1.5, 0.5)$ , the depth-1 decision tree policy cycles between taking an information gathering action  $x \leq 1$  and moving down to get a positive reward for which it gets the returns :

$$\begin{aligned}
 V^{\pi_{\mathcal{T}_1}}(s_g, o_0) &= \zeta + \gamma + \gamma^2 \zeta + \gamma^3 \dots \\
 &= \sum_{t=0}^{\infty} \gamma^{2t} \zeta + \sum_{t=0}^{\infty} \gamma^{2t+1} \\
 &= \frac{\zeta + \gamma}{1 - \gamma^2}
 \end{aligned}$$

At initialization, in either of the base states  $s_0 = (0.5, 0.5)$  and  $s_2 = (1.5, 1.5)$ , the value of the depth-1 decision tree policy is the return when taking one information gathering action  $x \leq 1$ , then moving right or down, then following the policy from the goal state  $s_g$  :

$$\begin{aligned}
 V^{\pi_{\mathcal{T}_1}}(s_0, o_0) &= \zeta + \gamma 0 + \gamma^2 V^{\pi_{\mathcal{T}_1}}(s_g, o_0) \\
 &= \zeta + \gamma^2 V^{\pi_{\mathcal{T}_1}}(s_g, o_0) \\
 &= V^{\pi_{\mathcal{T}_1}}(s_2, o_0)
 \end{aligned}$$

Similarly, the value of the best depth-1 decision tree policy in state  $s_1 = (0.5, 1.5)$  is the value of taking one information gathering action then moving right to  $s_2$  then following the policy in  $s_2$  :

$$\begin{aligned}
 V^{\pi_{T_1}}(s_1, o_0) &= \zeta + \gamma 0 + \gamma^2 V^{\pi_{T_1}}(s_2, o_0) \\
 &= \zeta + \gamma^2 V^{\pi_{T_1}}(s_2, o_0) \\
 &= \zeta + \gamma^2 (\zeta + \gamma^2 V^{\pi_{T_1}}(s_g, o_0)) \\
 &= \zeta + \gamma^2 \zeta + \gamma^4 V^{\pi_{T_1}}(s_g, o_0)
 \end{aligned}$$

Since the probability of being in any base states at initialization given that the agent observe  $o_0$  is the probability of being in any base states at initialization, we can write :

$$\begin{aligned}
 V^{\pi_{T_1}}(o_0) &= \frac{1}{4} V^{\pi_{T_1}}(s_g, o_0) + \frac{2}{4} V^{\pi_{T_1}}(s_2, o_0) + \frac{1}{4} V^{\pi_{T_1}}(s_1, o_0) \\
 &= \frac{1}{4} \frac{\zeta + \gamma}{1 - \gamma^2} + \frac{2}{4} (\zeta + \gamma^2 \frac{\zeta + \gamma}{1 - \gamma^2}) + \frac{1}{4} (\zeta + \gamma^2 \zeta + \gamma^4 \frac{\zeta + \gamma}{1 - \gamma^2}) \\
 &= \frac{1}{4} \frac{\zeta + \gamma}{1 - \gamma^2} + \frac{2}{4} (\frac{\zeta + \gamma^3}{1 - \gamma^2}) + \frac{1}{4} (\frac{\zeta + \gamma^5}{1 - \gamma^2}) \\
 &= \frac{4\zeta + \gamma + 2\gamma^3 + \gamma^5}{4(1 - \gamma^2)}
 \end{aligned}$$

□

We can now plot the POIBMDP objective values of the different policies corresponding to trees for the grid world MDP as functions of  $\zeta$  when we fix  $\gamma = 0.99$ . When  $\gamma = 0.99$ , despite objective values being very similar for the depth-1 and unbalanced depth-2 tree, we now know that **a depth-1 tree is the optimal deterministic partially observable POIBMDP policy for  $0 < \zeta < 1$ .**

Interestingly, two POMDP challenges described in (cite) can also be observed in Figure (cite). First, there is a whole range of  $\zeta$  values for which the stochastic policy is optimal. Second, for e.g.  $\zeta = 0.5$ , while the depth-1 tree is the optimal deterministic partially observable policy, the value of state  $(s_2, o_0) = (1.5, 1.5, 0, 2, 0, 2)$  is not maximized by this policy but by the sub-optimal policy that always goes down.

Now that we know to what values to set  $\zeta$  in our POIBMDPs, we can apply reinforcement learning algorithms to learn partially observable deterministic policies corresponding to the optimal depth-1 decision tree.



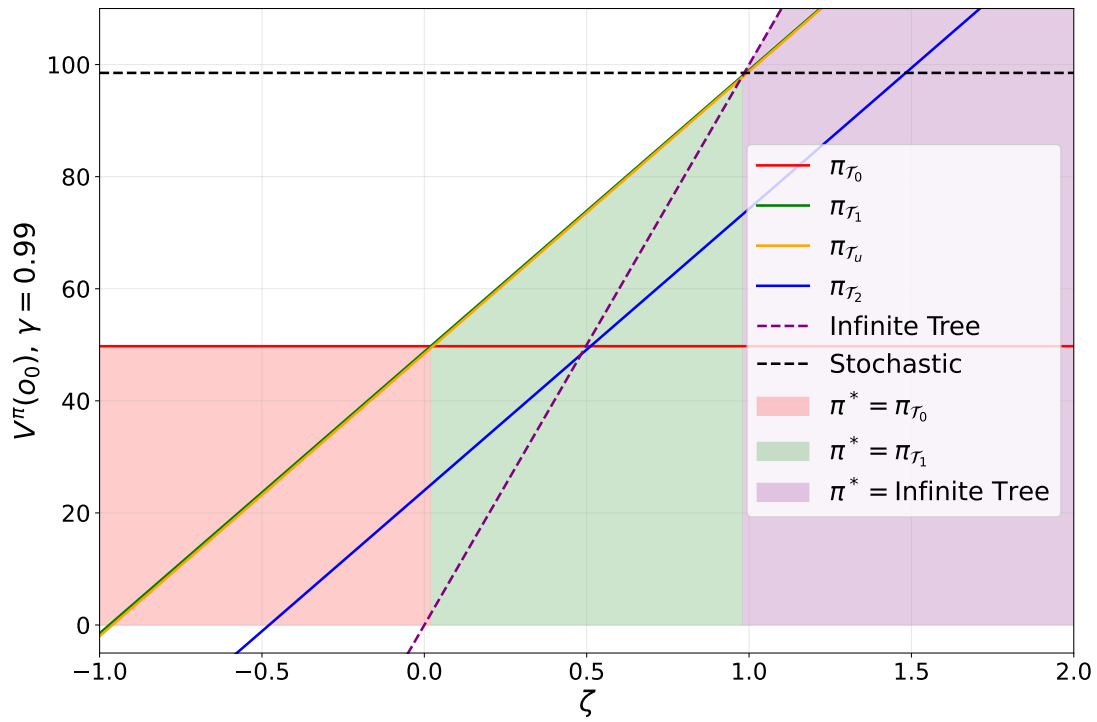


FIGURE 3.3 – POIBMDP objective values of different policies as functions of  $\zeta$ . Shaded areas show the optimal policies in different ranges of  $\zeta$  values.

## 3.2 Results

Unfortunately, our results are negative and show that reinforcement learning fails for the aforementioned problem. Let us understand why.

### 3.2.1 Experimental Setup

**Baselines :** we consider two groups of RL algorithms. The first group is standard tabular RL naively applied to POIBMDPs; Q-learning (cite)(cite), Sarsa (cite)(cite), and Policy Gradient (cite)(cite). (cite) and (cite) already studied the theoretical and practical implications of applying standard RL directly to POMDP by considering that partial observations are full MDP states. (cite) proved that Q-Learning will converge but without optimality guarantees. (cite) showed empirically that Sarsa- $\lambda$  (cite), a version of Sarsa with some sort of memory, can learn good deterministic solutions to POMDP which makes it a good candidate for our problem.

We also use a vanilla tabular Policy Gradient baseline (cite) with softmax policy, which to the best of our knowledge, nobody studied in our setting. In theory the Policy Gradient algorithm should not be a good candidate for our problem since it searches for stochastic policies that we showed can be better than our seeked depth-1 decision tree policy (c.f. Figure (cite)).

More recently, RL algorithms were developed specially for learning policies in POMDPs. Those algorithms are called asymmetric RL which authors of the original IBMDP paper (cite) use without knowing. Asymmetric algorithms leverage the fact that, even if the deployed policy should depend only on partial observation, nothing forbids the use of full state information during training if the latter is available. Asymmetric Q-learning (cite) makes use of full-state-action value function  $U : S \times O \rightarrow \mathbb{R}$  to use as a temporal difference error target (cite) when updating the  $Q : O \times A \rightarrow \mathbb{R}$  of interest. In Algorithm (cite), we write the Asymmetric Q-learning algorithm.

In addition to the traditional tabular RL algorithms, we also apply Asymmetric Q-learning and Asymmetric Sarsa and JSJ, the tabular algorithm from (cite)(cite) which is equivalent to a tabular Asymmetric Policy Gradient algorithm.

We use at least 200 000 time steps to train each agent. Each agent is trained on 100 seeds on each POIBMDP.

**Hyperparameters :** For all baselines we use, when applicable, exploration rates  $\epsilon = 0.3$  and learning rates  $\alpha = 0.1$ .

**Algorithm 2 : Asymmetric Q-Learning (cite)**


---

**Data :** POMDP  $\mathcal{M}_{po} = \langle X, O, A, R, T, T_0, \Omega \rangle$ , learning rates  $\alpha_u, \alpha_q$ , exploration rate  $\epsilon$

**Result :**  $\pi : O \rightarrow A$

Initialize  $U(x, a) = 0$  for all  $x \in X, a \in A$

Initialize  $Q(o, a) = 0$  for all  $o \in O, a \in A$

**for each episode do**

    Initialize state  $x_0 \sim T_0$

    Initialize observation  $o_0 \sim \Omega(x_0)$

**for each step  $t$  do**

        Choose action  $a_t$  using  $\epsilon$ -greedy :  $a_t = \arg \max_a Q(o_t, a)$  with prob.  $1 - \epsilon$

        Take action  $a_t$ , observe  $r_t = R(x_t, a_t)$ ,  $x_{t+1} \sim T(x_t, a_t)$ , and  $o_{t+1} \sim \Omega(x_{t+1})$

$y \leftarrow r + \gamma U(x_{t+1}, a_t) - U(x_t, a_t)$  // TD target

$U(x_t, a_t) \leftarrow (1 - \alpha_u)U(x_t, a_t) + \alpha_u y$

$Q(o_t, a_t) \leftarrow (1 - \alpha_q)Q(o_t, a_t) + \alpha_q y$

$x_t \leftarrow x_{t+1}$

$o_t \leftarrow o_{t+1}$

**end**

**end**

$\pi(o) = \arg \max_a Q(o, a)$  // Extract greedy policy

---

**Metrics :** we will consider two metrics. First, the sub-optimality gap of the learned policy POIBMDP value with respect to the optimal deterministic partially observable POIBMDP policy :  $|V^{\pi^*}(o_0) - V^\pi(o_0)|$  Because we know the whole POIBMDP model that we can represent exactly as tables; and because we know for each  $\zeta$  value the POIBMDP objective value of the optimal partially observable policy (c.f. Figure); we can report the *exact* sub-optimality gaps.

Second, we consider the distribution of the learned trees structure over the 100 training seeds. Indeed, since for every POIBMDP we train each baseline 100 times, we obtain 100 partially observable deterministic policies from which we can extract the equivalent 100 decision tree policies using Algorithm (cite). This helps understand which trees RL tend to learn.

### 3.2.2 Can RL baselines solve POIBMDPs?

In Figure (cite), we plot the sub-optimality gaps—averaged over 100 seeds—of learned policies during training. We do so for 200 different POIBMDPs where we change the reward for information gathering actions : we sample  $\zeta$  uniformly in  $[-1, 2]$ . For (asymmetric) Q-learning and (asymmetric) Sarsa, we extract the greedy policies from

---

**Algorithm 3 :** JSJ algorithm. Uses Monte Carlo estimates of the average reward value functions to perform policy improvements (cite)

---

**Data :** POMDP  $\mathcal{M}_{po} = \langle X, O, A, R, T, T_0, \Omega \rangle$ , learning rate  $\alpha$ , policy parameters  $\theta$ , number of trajectories  $N$

**Result :** Stochastic partially observable policy  $\pi_\theta : O \rightarrow \Delta A$

Initialize policy parameters  $\theta$

Initialize  $Q(o, a) = 0$  for all observations  $o$  and actions  $a$

**for each episode do**

**for**  $i = 1$  to  $N$  **do**

    Generate trajectory  $\tau_i = (s_0, a_0, r_0, s_1, a_1, r_1, \dots, s_T)$  following  $\pi_\theta$

**for each timestep**  $t$  **in trajectory**  $\tau_i$  **do**

$G_t \leftarrow \sum_{k=t}^T \gamma^{k-t} r_k$  // Compute return

      Store  $(o_t, a_t, G_t)$  for later averaging

**end**

**end**

**for each unique observation-action pair**  $(o, a)$  **do**

$Q(o, a) \leftarrow \frac{1}{|\{(o, a)\}|} \sum_{(o, a, G)} G$  // Monte Carlo estimate

**end**

**for each observation**  $o$  **do**

**for each action**  $a$  **do**

$\pi_1(a|o) \leftarrow 1.0$  if  $a = a'$   $Q(o, a')$ ,  $0.0$  otherwise // Deterministic

      policy from  $Q$ -values

$\pi(a|o) \leftarrow (1 - \alpha)\pi(a|o) + \alpha\pi_1(a|o)$  // Policy improvement step

**end**

**end**

  Reset  $Q(o, a) = 0$  for all observations  $o$  and actions  $a$  // Reset for next episode

**end**

---

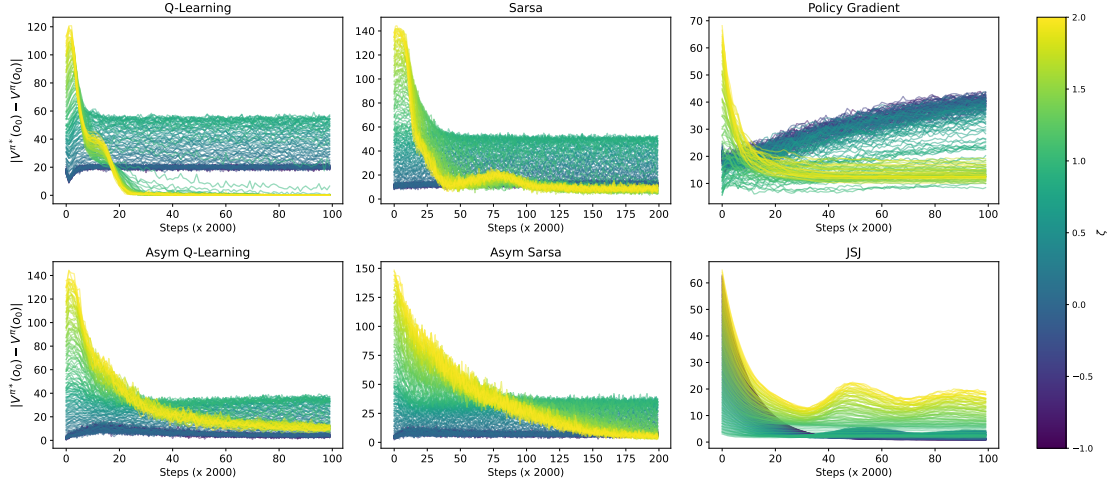


FIGURE 3.4 – Baselines learning curves of partially observable policies in POIBMDPs. In each subplot, each single learning curve is colored by the value of  $\zeta$  in the corresponding POIBMDP in which learning occurs. Each single learning curve represent the sub-optimality gap averaged over 100 seeds. So for each baseline we ran a total of  $200 \times 100$  single training run.

the learned Q-functions and evaluate them every 2000 steps. For the Policy Gradient and the JSJ baselines we directly evaluate the learned stochastic policies.

In Figure (cite), we plot the distributions over the final learned trees in function of  $\zeta$  from the above runs. For the Policy Gradient and JSJ baselines, we extract the tree from the policy that is greedy w.r.t to the action probabilities as Algorithm 6 requires a deterministic partially observable policy to return a tree policy.

We observe that, despite all runs converging, independently of the  $\zeta$  values, not all runs fully minimize the sub-optimality gap. In particular all RL algorithms seem to consistently minimize the gap, i.e. learn the optimal policy or Q-function, for  $\zeta \in [-1, 0]$ , where the optimal policy is the depth-0 tree. For  $\zeta \in [1, 2]$ , where the optimal policy is to repeat taking information gathering actions, only the Q-learning baseline learns the optimal policy. In  $\zeta \in ]0, 1[$ , where the depth-1 tree is optimal, no baseline can consistently learn the optimal solution. However, we observe that asymmetric versions of Q-learning and Sarsa have found the optimal policy more frequently than other baselines. One interpretation of this phenomenon is that the learning in POIBMDPs is very difficult and so agents tend to converge to trivial policies, e.g., repeating the same base action.

Next, we quantify how difficult it is to do RL to learn policies in POIBMDPs.

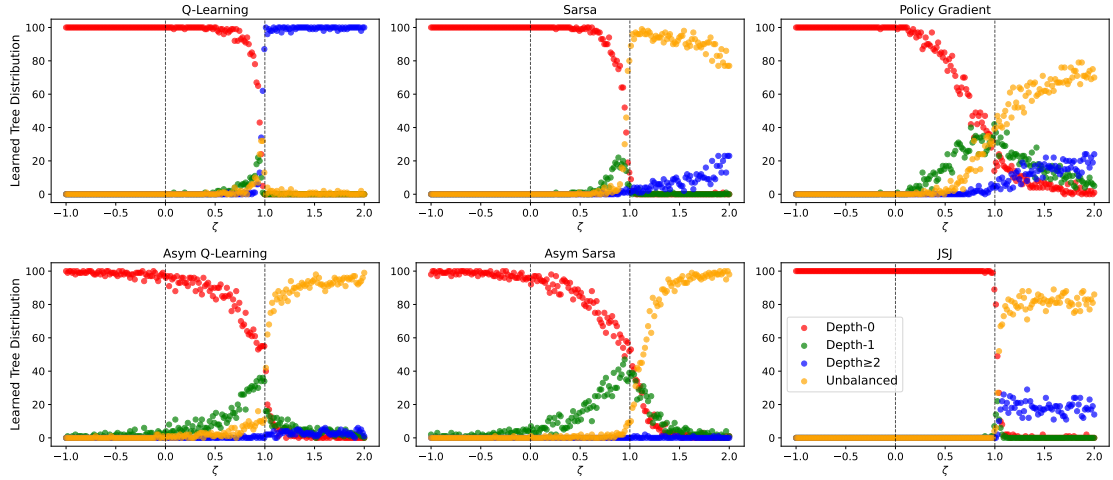


FIGURE 3.5 – Distributions of final tree policies learned across the 100 seeds. For each  $\zeta$  value, there are four color points each representing the share of depth-0 trees (red), depth-1 trees (green), unbalanced depth-2 trees (orange) and depth-2 trees (blue). Note that those trees do not necessarily correspond to trees in Figure (cite). We simply check the structure of trees and not if they are the best-in-class tree.

### 3.2.3 How difficult is it to learn in POIBMPDs?

As discussed in the previous chapter (cite), POIBMDPs inherit many difficulties from POMDPs. To prove that solving POIBMDPs, even though there are only a handful of states, observations, and actions (16, 9, and 6 respectively), we compare the success rate of RL baselines when applied to a POIBMDP with  $\gamma = 0.99$ ,  $\zeta = 0.5$  and when applied to the corresponding fully observable IBMDP. Because we know that for those values of  $\zeta$  and  $\gamma$ , the optimal deterministic partially observable policy is a depth-1 tree and the optimal fully observable policy is a tabular policy from Figure (cite); we can compute the success rate of an RL algorithm as the proportion of learned policies that match their optimal counterparts.

Furhtermore; we also consider advanced tricks for tabular RL in order to properly assess the difficulty of learning in POIBMDPs compared to learning in IBMDPs :

1. Optimistic Q-function (cite)
2. Entropy regularization (cite)
3. Eligibility traces (cite)
4.  $\epsilon$ -decay (cite)

In particular, for each of the aforementioned RL baselines; we learn policies using a wide variety of hyperparameters. For example, in Table (cite) we provide the detailed

hyperparameter space of Asymmetric Sarsa, which induce a total of 1152 instances of Asymmetric Sarsa, and in Table (cite) we provide the hyperparameter space sizes for all baselines We run each baseline using each hyperparamters combination 10 times on 1 million timesteps.

TABLEAU 3.1 – Asymmetric Sarsa Hyperparameter Space (768 combinations)

Hyperparameter	Values	Description
Epsilon Schedules	(0.3, 1), (0.3, 0.99), (1, 1)	Initial exploration and decrease rate
Epsilon Schedules	(0.1, 1), (0.1, 0.99), (0.3, 0.99)	Initial exploration and decrease rate
Lambda	0.0, 0.3, 0.6, 0.9	Eligibility trace decay
Learning Rate $U$	0.001, 0.005, 0.01, 0.1	$Q$ learning rate
Learning Rate $Q$	0.001, 0.005, 0.01, 0.1	$U$ learning rate
Optimistic	True, False	Optimistic initialization

TABLEAU 3.2 – Summary of RL baselines Hyperparameter Spaces

Algorithm	Problem	Total Hyperparameter Combinations
Policy Gradient	PO/IBMDP	420
JSJ	POIBMDP	15
Q-learning	PO/IBMDP	192
Asym Q-learning	POIBMDP	768
Sarsa	PO/IBMDP	192
Asym Sarsa	POIBMDP	768

Best hyperparams of Asymmetric Q-learning solves the POIBMDP 10/10 times :

$\epsilon_{start} : 1.0$

$\epsilon_{decay} : 0.99$

$batch\_size : 1$

$\lambda : 0.0$

$lr_o : 0.01$

$lr_v : 0.1$

$optimistic : False$

Asym Sarsa (10/10) :

$\epsilon_{start} : 1.0$

$\epsilon_{decay} : 0.99$

$batch\_size : 1$

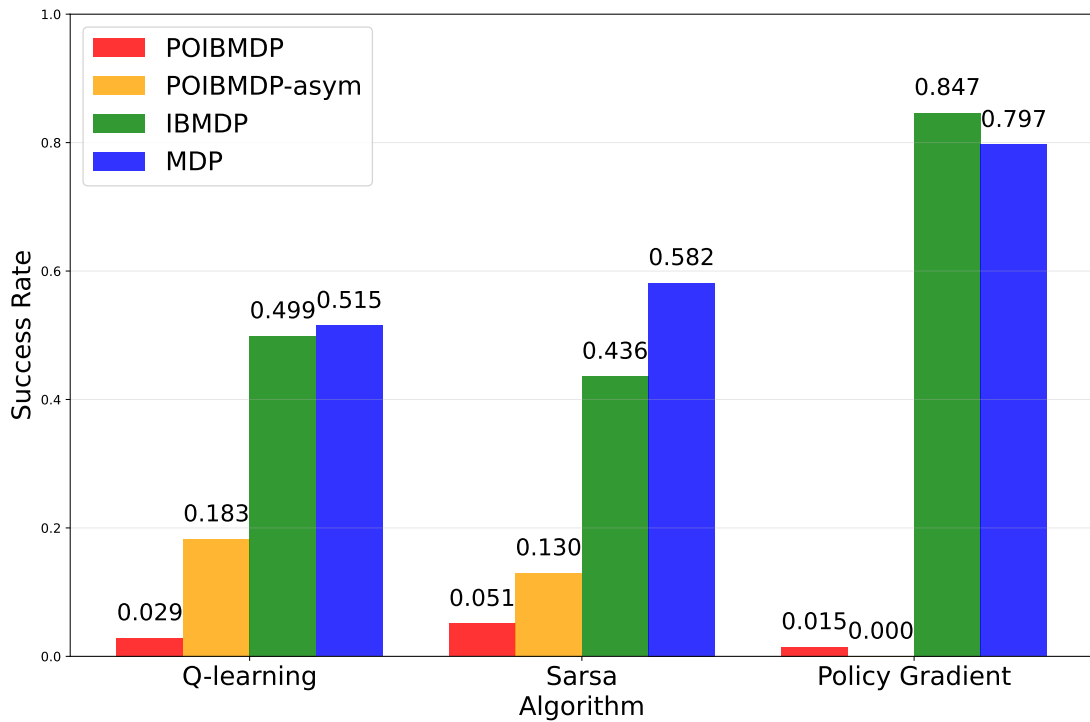


FIGURE 3.6 – Success rates of different RL algorithms over thousands of runs when applied to a POIBMDP and its fully observable IBMDP counterpart

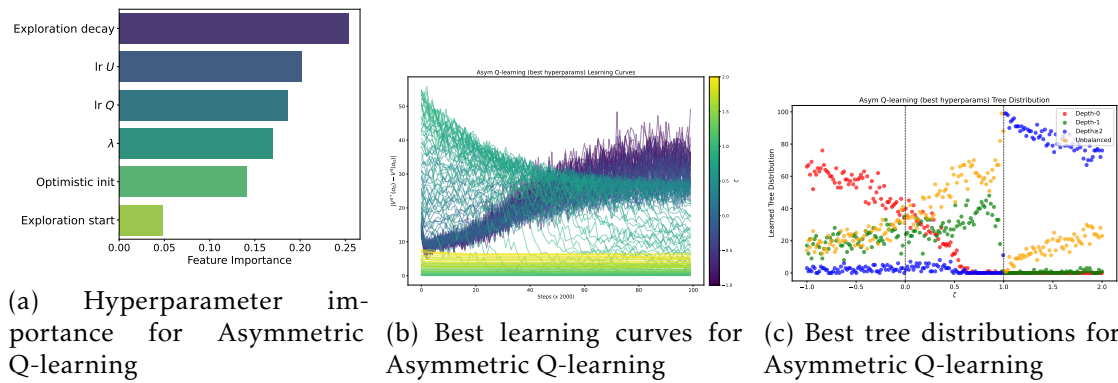


FIGURE 3.7 – Asymmetric Q-learning performance analysis : hyperparameter importance, learning curves, and tree distributions across different POIBMDP configurations.



```

lambda:0.0
lro : 0.1
lrv : 0.005
optimistic : False

```

Best hyperparameters for PG (4/10) :

```

lr : 0.05
tau : 0.1
eps : 0.1
nssteps : 2000

```

The key observations from Figure (cite) is that POIBMDPs are way harder than their IBMDPs counterparts. Even though asymmetry seems to increase performances; learning a decision tree policy for a simple grid world directly with RL using the framework of POIBMDP seem way to difficult and costly as successes might require a million steps for such a seemingly simple problem. An other difficulty in practice that we did not cover is the choice of information gathering actions. For the grid world MDP, choosing good IGAs ( $x \leq 1$  and  $y \leq 1$ ) is simple but what about more complicated MDPs? Even after choosing candidate IGAs; RL with large action space is known to be an already difficult problem.

### 3.3 Conclusion

We suggest researchers and practitioners interested in decision tree policies for MDPs to forget about the direct RL approach and focus on the indirect imitation-based approach (cite) that despite being prone to failures too (cite), is cheaper sample and computation wise as one can re-use pretrained expert policies.

If for some reason one has to avoid imitation; fitting parametric trees, trees whose depth and nodes are predefined and where only thresholds are optimized, then one can use SYMPOL (cite) or Soft decision trees (cite). The latter approach can be beneficial when the user has lots of domain knowledge on what nodes should be in a decision tree policy.

In the next chapter, we show that when the base MDP's transitions are independent of the action, then POIBMDPs are actually *fully* observable which can lead to new decision tree induction for supervised learning.



## When transitions are uniform POIBMDPs are fully observable

In this section we show that decision tree induction for classification problems can be formulated as solving POIBMDPs. Indeed, a supervised classification problem (cite) can be formulated as an MDP where actions are class labels and states are training data. The reward at every step is 1 if the correct label was predicted and 0 otherwise. In this MDP, the transitions are independent of the actions : the next state is given by uniformly sampling a new training datum. This implies that the resulting POIBMDP is actually *fully* observable and that traditional RL algorithms should perform well.

**Definition 5** (Classification Markov Decision Process). *Given a set of  $N$  examples denoted  $\mathcal{E} = \{(x_i, y_i)\}_{i=1}^N$  where each datum  $x_i$  is described by a set of  $p$  features and  $y_i \in \mathbb{Z}^m$  is the label associated with  $x_i$ , a Classification Markov Decision Process is an MDP (cite)  $\langle S, A, R, T, T_0 \rangle$  where :*

- *the state space is  $S = \{x_i\}_{i=1}^N$ , the set of data features*
- *the action space is  $A = \mathbb{Z}^m$ , the set of unique labels*
- *the reward function is  $R : S \times A \rightarrow \{0, 1\}$  with  $R(s = x_i, a) = 1_{a=y_i}$*
- *the transition function is  $T : S \times A \rightarrow \Delta S$  with  $T(s, a, s') = \frac{1}{N} \quad \forall s, a, s'$*
- *the initial distribution is  $T_0(s_0 = s) = \frac{1}{N}$*

It is easy to see that policies that maximize the expected discounted cumulative reward objective (cite) are classifiers that maximize the prediction accuracy because  $\sum_{i=1}^N 1_{\pi(x_i)=y_i} = \sum_{i=1}^N R(x_i, \pi(x_i))$ . We defer the formal proof in the next part of the manuscript in which we extensively study supervised learning problems. Learning a classifier with a reward signal can also be done with contextual bandits (cite).

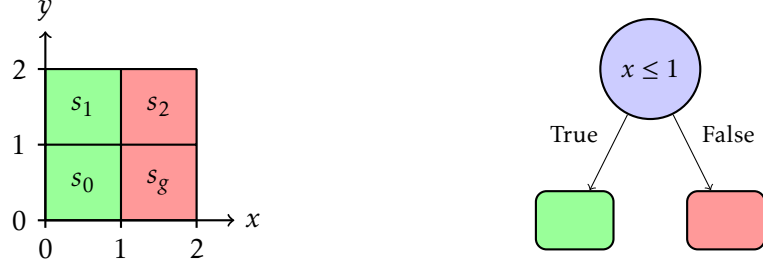


FIGURE 4.1 – Classification MDP optimal actions. In this Classification MDP, there are four data to assign either a green or red label. On the right, there is the unique optimal depth-1 tree for this particular Classification MDP. This depth-1 tree also maximizes the accuracy on the corresponding classification tasks.

To learn a decision tree *classifier*, we use the POIBMDP formalism and show, in the particular case where the base MDP is a Classification MDP (cite), that associated POIBMDPs are in fact MDPs with stochastic transitions.

**Definition 6** (Classification POIBMDP). *Given a Classification MDP (cite)  $\langle \{x_i\}_{i=1}^N, \mathbb{Z}^m, R, T, T_0 \rangle$ , and an associated POIBMDP (cite)  $\langle S, O, A, A_{info}, R, \zeta, T_{info}, T, T_0 \rangle$ , a Classification POIBMDP is an MDP :*

$$\left\langle \underbrace{\overbrace{O}}_{\text{State space}}, \underbrace{\mathbb{Z}^m, A_{info}}_{\text{Action space}}, \underbrace{\overbrace{R, \zeta}}_{\text{Reward function}}, \underbrace{\overbrace{\mathcal{P}, \mathcal{P}_0}}_{\text{Transition kernels}} \right\rangle$$

- $O$  is the set of possible observations in  $[L_1, U_1] \times \dots \times [L_d, U_d] \times [L_1, U_1] \text{ times } \dots \times [L_d, U_d]$  where  $L_j$  is the minimum value of feature  $j$  over all data  $x_i$  and  $U_j$  the maximum
- $\mathbb{Z}^m \cup A_{info}$  is action space : actions can be label assignments in  $\mathbb{Z}^m$  or bounds refinement in  $A_{info}$
- The reward for assigning label  $a \in \mathbb{Z}^m$  when observing some observation  $o = (L'_1, U'_1, \dots, L'_d, U'_d)$  is the proportion of training data satisfying the bounds and having label  $a$  :  $R(o, a) = \frac{|\{x_i : L'_j \leq x_{ij} \leq U'_j \forall i, j\} \cap \{x_i : y_i = a \forall i\}|}{|\{x_i : L'_j \leq x_{ij} \leq U'_j \forall i, j\}|}$ . The reward for taking an information gathering action that refines bounds is  $\zeta$
- The transition kernel is  $\mathcal{P} : O \times (\mathbb{Z}^m \cup A_{info}) \rightarrow \Delta O$  where :
  - For  $a \in \mathbb{Z}^m$  :  $\mathcal{P}(o, a, (L_1, U_1, \dots, L_d, U_d)) = 1$  (reset to full bounds)
  - For  $a = (k, v) \in A_{info}$  : from  $o = (L'_1, U'_1, \dots, L'_d, U'_d)$ , the MDP will transit to  $o_{left} = (L'_1, U'_1, \dots, L_k, v, \text{dots}, L'_d, U'_d)$  (resp.  $o_{right} = (L'_1, U'_1, \dots, U_k, v, \text{dots}, L'_d, U'_d)$ ) with

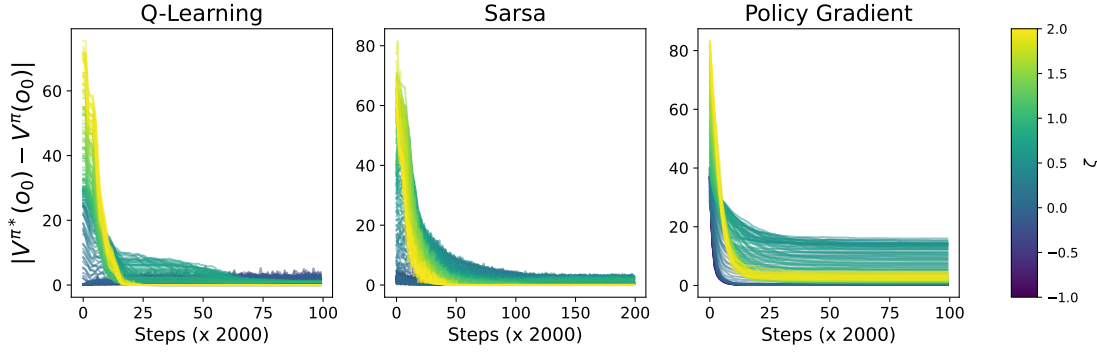


FIGURE 4.2

$$probability \frac{|\{x_i: L'_j \leq x_{ij} \leq U'_j \forall j \wedge x_{ik} \leq v\}|}{|\{x_i: L'_j \leq x_{ij} \leq U'_j \forall j\}|} \text{ (resp. } \frac{|\{x_i: L'_j \leq x_{ij} \leq U'_j \forall j \wedge x_{ik} > v\}|}{|\{x_i: L'_j \leq x_{ij} \leq U'_j \forall j\}|} \text{)}$$

Now we can learn a decision tree *classifier* for a supervised learning task as follows. Given a set of examples and labels, we can formulate the associated Classification MDP (cite). Then we can provide some information gathering actions and a reward  $\zeta$  for those to define a Classification POIBMDP (cite). Finally, we can learn a deterministic policy  $\pi : O \rightarrow A \cup A_{info}$  to obtain a decision tree classifier. Since we have shown that Classification POIBMDPs are MDPs, we expect that RL baselines should perform well to solve Classification POIBMDPs. We expect that, compared to general POIBMDPs studied in the previous chapter where some information was hidden to the agent, RL baselines will consistently learn the optimal decision tree classifiers.

## 4.1 How well can RL baselines learn in Classification POIBMDPs?

Similarly to the previous chapter, we are interested in a very simple Classification POIBMDP that corresponds to building a tree for the supervised problem :

$$\mathcal{X} = \{(0.5, 0.5), (0.5, 1.5), (1.5, 1.5), (1.5, 0.5)\}$$

$$y = \{0, 0, 1, 1\}$$

We illustrate the associated Classification MDP in Figure (cite). We construct Classification POIBMDPs with  $\gamma = 0.99$ , 200 values of  $\zeta \in [0, 1]$  and IGAs  $x \leq 1$  and  $y \leq 1$ . Since Classification POIBMDPs are MDPs, we do not need to analyze asymmetric and JSJ (cite) baselines.

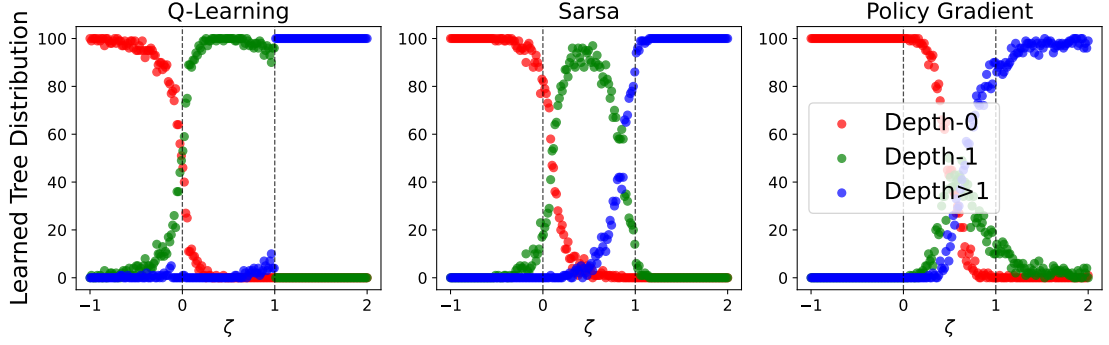


FIGURE 4.3

Fortunately this time, compared to general POIBMDPs, RL can be used to retrieve optimal policies in Classification POIBMDPs equivalent to decision tree classifiers. We observe on Figure (cite) that both Q-learning and Sarsa consistently minimize the sub-optimality gap. Hence they are able to retrieve the optimal depth-1 decision tree classifier from Figure (cite).

In this part of the manuscript, we have highlighted the challenges of reinforcement learning of decision tree policies for MDPs. We showed that this task was particularly challenging, even when an optimal decision tree policy exists. We saw that RL algorithms are unable to find such optimal trees, even for very simple problems such as a  $2 \times 2$  grid world. However, we also showed that for some MDPs, such as ones that encode a supervised learning task, RL algorithms could retrieve optimal decision tree classifiers. This was shown only on a very simple classification task with four data and two classes. Furthermore, good candidate nodes—information gathering actions—were provided to the learning algorithms.

Now remains the question of scalability and competitiveness of such approaches to supervised learning for which strong baselines have existed for decades. In the next part of this thesis; we present a new decision tree induction that relies on solving MDPs exactly and choosing information gathering actions adaptively.

**Deuxième partie**

**An easier problem : Learning  
Decision Trees for MDPs that are  
Classification tasks**





## DPDT-intro

In supervised learning, decision trees are valued for their interpretability and performance. While greedy decision tree algorithms like CART remain widely used due to their computational efficiency, they often produce sub-optimal solutions with respect to a regularized training loss. Conversely, optimal decision tree methods can find better solutions but are computationally intensive and typically limited to shallow trees or binary features. We present Dynamic Programming Decision Trees (DPDT), a framework that bridges the gap between greedy and optimal approaches. DPDT relies on a Markov Decision Process formulation combined with heuristic split generation to construct near-optimal decision trees with significantly reduced computational complexity. Our approach dynamically limits the set of admissible splits at each node while directly optimizing the tree regularized training loss. Theoretical analysis demonstrates that DPDT can minimize regularized training losses at least as well as CART. Our empirical study shows on multiple datasets that DPDT achieves near-optimal loss with orders of magnitude fewer operations than existing optimal solvers. More importantly, ex-

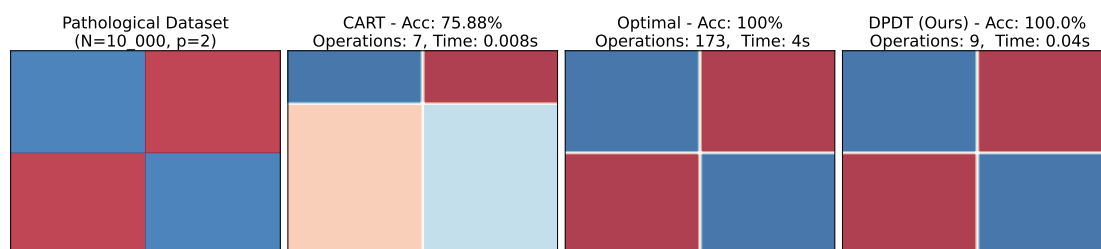


FIGURE 5.1 – Pathological dataset and learned depth-2 trees with their scores, complexities, runtimes, and decision boundaries.

tensive benchmarking suggests statistically significant improvements of DPDT over both CART and optimal decision trees in terms of generalization to unseen data. We demonstrate DPDT practicality through applications to boosting, where it consistently outperforms baselines. Our framework provides a promising direction for developing efficient, near-optimal decision tree algorithms that scale to practical applications.

## 5.1 Introduction

Decision trees [83, 82, 13] are at the core of various machine learning applications. Ensembles of decision trees such as tree boosting [39, 38, 20, 79] are the state-of-the-art for supervised learning on tabular data [44]. Human users can make sense of decision trees predictions [31, 62, 15] which allows for real-world applications when safety or trust is critical [89]. More recently, decision trees have been used to model sequential decision policies with imitation learning [kohler2024interpretable, 6] or directly with reinforcement learning (RL) [marton2024sympolsymbolictreebasedonpolicy, 97, 107, 106].

To motivate the design of new decision tree induction algorithms, Figure 5.1 exhibits a dataset for which existing greedy algorithms are suboptimal, and optimal algorithms are computationally expensive. The dataset is made up of  $N = 10^4$  samples in  $p = 2$  dimensions that can be perfectly labeled with a decision tree of depth 2. When running CART [13], greedily choosing the root node yields a suboptimal tree. This is because greedy algorithms compute locally optimal splits in terms of information gain. In our example, the greedy splits always give two children datasets which themselves need depth 2 trees to be perfectly split. On the other hand, to find the root node, an optimal algorithm such as [68] iterates over all possible splits, that is,  $N \times p = 20,000$  operations to find one node of the solution tree.

In this work, we present a framework for designing non-greedy decision tree induction algorithms that optimize a regularized training loss nearly as well as optimal methods. This is achieved with orders of magnitude less operations, and hence dramatic computation savings. We call this framework “Dynamic Programming Decision Trees” (DPDT). For every node, DPDT heuristically and dynamically limits the set of admissible splits to a few good candidates. Then, DPDT optimizes the regularized training loss with some depth constraints. Theoretically, we show that DPDT minimizes the empirical risk at least as well as CART. Empirically, we show that on all tested datasets, DPDT can reach 99% of the optimal regularized train accuracy while using thousands times less operations than current optimal solvers. More importantly, we follow [44] methodology

to benchmark DPDT against both CART and optimal trees on hard datasets. Following the same methodology, we compare boosted DPDT [37] to boosted CART and to some deep learning methods and show clear superiority of DPDT.

## 5.2 Related Work

To learn decision trees, greedy approaches like CART [13] iteratively partition the training dataset by taking splits optimizing a local objective such as the Gini impurity or the entropy. This makes CART suboptimal with respect to training losses [74]. But CART remains the default decision tree algorithm in many machine learning libraries such as [77, 20, 52, 113] because it can scale to very deep trees and is very fast. To avoid overfitting, greedy trees are learned with a maximal depth or pruned a posteriori [13, chapter 3]. In recent years, more complex optimal decision tree induction algorithms have shown consistent gains over CART in terms of generalization capabilities [9, 104, 28].

Optimal decision tree approaches optimize a regularized training loss while using a minimal number of splits [9, 3, 105, 68, 28, 29, 60, 19]. However, direct optimization is not a convenient approach, as finding the optimal tree is known to be NP-Hard [50]. Despite the large number of algorithmic tricks to make optimal decision tree solvers efficient [28, 68], their complexity scales with the number of samples and the maximum depth constraint. Furthermore, optimal decision tree induction algorithms are usually constrained to binary-features dataset while CART can deal with any type of feature. When optimal decision tree algorithms deal with continuous features, they can usually learn only shallow trees, e.g. Quant-BnB [68] can only compute optimal trees up to depth 3. PySTreeD, the latest optimal decision tree library [60], can compute decision trees with depths larger than three but uses heuristics to binarize a dataset with continuous features during a pre-processing step. Despite their limitations to binary features and their huge computational complexities, encouraging practical results for optimal trees have been obtained [73, 59, 22, 61]. Among others, they show that optimal methods under the same depth constraint (up to depth four) find trees with 1–2% greater test accuracy than greedy methods.

In this work, we only consider the induction of nonparametric binary depth-constrained axis-aligned trees. By nonparametric trees, we mean that we only consider tree induction algorithms that optimize both features and threshold values in internal nodes of the tree. This is different from the line of work on Tree Alternating Optimization (TAO) algorithm [17, 112, 18] that only optimizes tree nodes threshold values for fixed nodes

features similarly to optimizing neural network weights with gradient-based methods.

There exist many other areas of decision tree research [63] such as inducing non-axis parallel decision trees [75, 51], splitting criteria of greedy trees [61], different optimization of parametric trees [76, 79], or pruning methods [34, 72].

Our work is not the first to formulate the decision tree induction problem as solving a Markov decision process (MDP) [32, 40, 97, 19]. Those works formulate tree induction as solving a partially observable MDP and use approximate algorithms such as Q-learning [40] or deep Q-learning [97] to solve them in an online fashion one datum from the dataset at a time. In a nutshell, our work, DPDT that we present next, is different in that it builds a stochastic and fully observable MDP that can be explicitly solved with dynamic programming. This makes it possible to solve exactly the decision tree induction problem.

## DPDT-paper

### 6.1 Decision Trees for Supervised Learning

Let us briefly introduce some notations for the supervised classification problem considered in this paper. We assume that we have access to a set of  $N$  examples denoted  $\mathcal{E} = \{(x_i, y_i)\}_{i=1}^N$ . Each datum  $x_i$  is described by a set of  $p$  features.  $y_i \in \mathcal{Y}$  is the label associated with  $x_i$ .

A decision tree is made of two types of nodes : split nodes that are traversed, and leaf nodes that finally assign a label. To predict the label of a datum  $x$ , a decision tree  $T$  sequentially applies a series of splits before assigning it a label  $T(x) \in \mathcal{Y}$ . In this paper, we focus on binary decision trees with axis-aligned splits as in [13], where each split compares the value of one feature with a threshold.

Our goal is to learn a tree that generalizes well to unseen data. To avoid overfitting, we constrain the maximum depth  $D$  of the tree, where  $D$  is the maximum number of splits that can be applied to classify a data. We let  $\mathcal{T}_D$  be the set of all binary decision trees of depth  $\leq D$ . Given a loss function  $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}_+$ , we induce trees with a regularized training loss defined by :

$$\begin{aligned}
T^* &= \operatorname{argmin}_{T \in \mathcal{T}_D} \mathcal{L}_\alpha(T), \\
T^* &= \operatorname{argmin}_{T \in \mathcal{T}_D} \frac{1}{N} \sum_{i=1}^N \ell(y_i, T(x_i)) + \alpha C(T),
\end{aligned} \tag{6.1}$$

where  $C : \mathcal{T} \rightarrow \mathbb{R}$  is a complexity penalty that helps prevent or reduce overfitting

such as the number of nodes [13, 68], or the expected number of splits to label a data[73]. The complexity penalty is weighted by  $\alpha \in [0, 1]$ . For supervised classification problems, we use the 0–1 loss :  $\ell(y_i, T(x_i)) = 1_{\{y_i \neq T(x_i)\}}$ . Please note while we focus on supervised classification problems in this paper, our framework extends naturally to regression problems.

We now formulate the decision tree induction problem 6.1 as finding the optimal policy in an MDP.

## 6.2 Decision Tree Induction as an MDP

Given a set of examples  $\mathcal{E}$ , the induction of a decision tree is made of a sequence of decisions : at each node, we must decide whether it is better to split (a subset of)  $\mathcal{E}$ , or to create a leaf node.

This sequential decision-making process corresponds to a Markov Decision Problem (MDP) [80]  $\mathcal{M} = \langle S, A, R_\alpha, P, D \rangle$ . A state is a pair made of a subset of examples  $X \subseteq \mathcal{E}$  and a depth  $d$ . Then, the set of states is  $S = \{(X, d) \in P(\mathcal{E}) \times \{0, \dots, D\}\}$  where  $P(\mathcal{E})$  denotes the power set of  $\mathcal{E}$ .  $d \in \{0, \dots, D\}$  is the current depth in the tree. An action  $A$  consists in creating either a split node, or a leaf node (label assignment). We denote the set of candidate split nodes  $\mathcal{F}$ . A split node in  $\mathcal{F}$  is a pair made of one feature  $i$  and a threshold value  $x_{ij} \in \mathcal{E}$ . So, we can write  $A = \mathcal{F} \cup \{1, \dots, K\}$ . From state  $s = (X, d)$  and a splitting action  $a \in \mathcal{F}$ , the transition function  $P$  moves to the next state  $s_l = (X_l, d + 1)$  with probability  $p_l = \frac{|X_l|}{|X|}$  where  $X_l = \{(x_i, y_i) \in X : x_i \leq x_{ij}\}$ , or to state  $s_r = (X \setminus X_l, d + 1)$  with probability  $1 - p_l$ . For a class assignment action  $a \in \{1, \dots, K\}$ , the chain reaches an absorbing terminal state with probability 1. The reward function  $R_\alpha : S \times A \rightarrow \mathbb{R}$  returns  $-\alpha$  for splitting actions and the proportion of misclassified examples of  $X - \frac{1}{|X|} \sum_{(x_i, y_i) \in X} \ell(y_i, a)$  for class assignment actions.  $\alpha \in [0, 1]$  controls the accuracy-complexity trade-off defined in the regularized training objective 6.1. The horizon  $D$  limits tree depth to  $D$  by forbidding class assignments after  $D$  MDP transitions.

The solution to this MDP is a deterministic policy  $\pi : S \rightarrow A$  that maximizes  $J_\alpha(\pi) = \mathbb{E}\left[\sum_{t=0}^D R_\alpha(s_t, \pi(s_t))\right]$ , the expected sum of rewards where the expectation is taken over transitions  $s_{t+1} \sim P(s_t, \pi(s_t))$  starting from initial state  $s_0 = (\mathcal{E}, 0)$ . Any such policy can be converted into a binary decision tree through a recursive extraction function  $E(\pi, s)$  that returns, either a leaf node with class  $\pi(s)$  if  $\pi(s)$  is a class assignment, or a tree with root node containing split  $\pi(s)$  and left/right sub-trees  $E(\pi, s_l)/E(\pi, s_r)$  if  $\pi(s)$  is a split. The final decision tree  $T$  is obtained by calling  $E(\pi, s_0)$  on the initial state  $s_0$ .

**Proposition 2** (Objective Equivalence). *Let  $\pi$  be a deterministic policy of the MDP and  $\pi^*$  be an optimal deterministic policy. Then  $J_\alpha(\pi) = -\mathcal{L}_\alpha(E(\pi, s_0))$  and  $T^* = E(\pi^*, s_0)$  where  $T^*$  is a tree that optimizes Eq. 6.1.*

This proposition is key as it states that the return of any policy of the MDP defined above is equal to the regularized training accuracy of the tree extracted from this policy. A consequence of this proposition is that when all possible splits are considered, the optimal policy will generate the optimal tree in the sense defined by Eq. (6.1). The proof is given in the Appendix 6.6.

## 6.3 Algorithm

We now present the Dynamic Programming Decision Tree (DPDT) induction algorithm. The algorithm consists of two essential steps. The first and most computationally expensive step constructs the MDP presented in Section 6.2. The second step solves it to obtain a policy that maximizes Eq.6.2 and that is equivalent to a decision tree. Both steps are now detailed.

### 6.3.1 Constructing the MDP

An algorithm constructing the MDP of section 6.2 essentially computes the set of all possible decision trees of maximum depth  $D$  which decision nodes are in  $\mathcal{F}$ . The transition function of this specific MDP is a directed acyclic graph. Each node of this graph corresponds to a state for which one computes the transition and reward functions. Considering all possible splits in  $\mathcal{F}$  does not scale. We thus introduce a state-dependent action space  $A_s$ , much smaller than  $A$  and populated by a splits generating function. In Figure ??, we illustrate the MDP constructed for the classification of a toy dataset using some arbitrary splitting function.

### 6.3.2 Heuristic splits generating functions

A split generating function is any function  $\phi$  that maps an MDP state, i.e., a subset of training examples, to a split node. It has the form  $\phi : S \rightarrow P(\mathcal{F})$ , where  $P(\mathcal{F})$  is the power set of all possible split nodes in  $\mathcal{F}$ . For a state  $s \in S$ , the state-dependent action space is defined by  $A_s = \phi(s) \cup \{1, \dots, K\}$ .

When the split generating function does not return all the possible candidate split nodes given a state, solving the MDP with state-dependent actions  $A_s$  is not guaranteed to yield the minimizing tree of Eq. 6.1, as the optimization is then performed on the

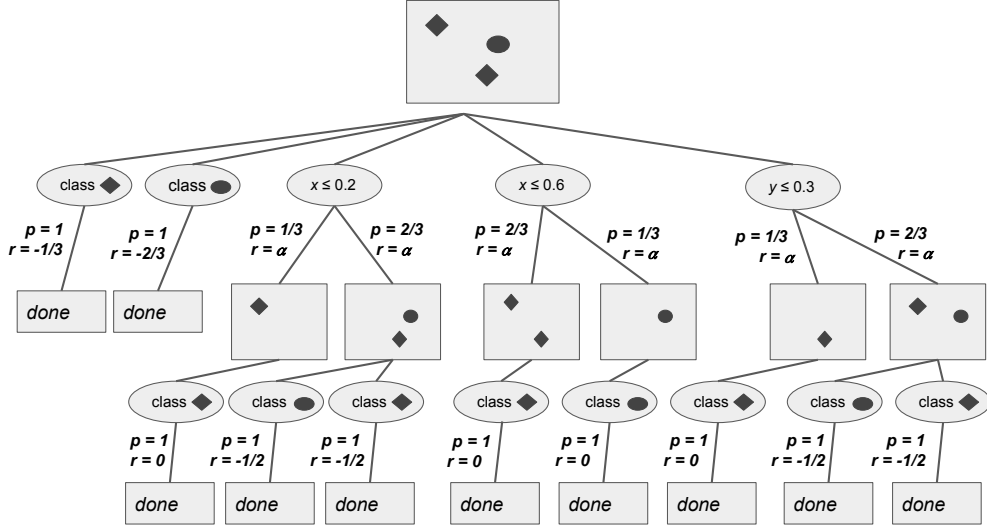


FIGURE 6.1 – Schematics of the MDP to learn a decision tree of depth 2 to classify a toy dataset with three samples, two features (x,y), and two classes (oval, diamond) and using an arbitrary splits generating function.

subset of trees of depth smaller or equal to  $D$ ,  $\mathcal{T}_D$ . We now define some interesting split generating functions and provide the time complexity of the associated decision tree algorithms. The time complexity is given in big-O of the number of candidate split nodes considered during computations.

**Exhaustive function.** When  $\mathcal{F} \subseteq \phi(s), \forall s \in S$ , the MDP contains all possible splits of a certain set of examples. In this case, *the optimal MDP policy is the optimal decision tree of depth at most  $D$* , and the number of states of the MDP would be  $O((2Np)^D)$ . Solving the MDP for  $A_s = \phi(s)$  is equivalent to running one of the optimal tree induction algorithms [104, 9, 60, 68, 105, 28, 3, 29, 59, 19]

**Top  $B$  most informative splits.** [11] proposed to generate splits with a function that returns, for any state  $s = (X, d)$ , the  $B$  most informative splits over  $X$  with respect to some information gain measure such as the entropy or the Gini impurity. The number of states in the MDP would be  $O((2B)^D)$ . When  $B = 1$ , *the optimal policy of the MDP is the greedy tree*. In practice, we noticed that the returned set of splits lacked diversity and often consists of splits on the same feature with minor changes to the threshold value.



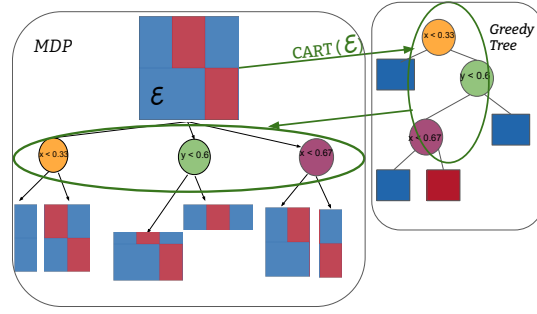


FIGURE 6.2 – How CART is used in DPDT to generate candidate splits given the example data in the current state.

**Calls to CART** Instead of returning the most informative split at each state  $s = (X, d)$ , we propose to find the most discriminative split, i.e. the feature splits that best predicts the class of data in  $X$ . We can do this by considering the split nodes of the greedy tree. In practice, we run CART on  $s$  and use the returned nodes as  $\phi(s)$ . We control the number of MDP states by constraining CART trees with a maximum number of nodes  $B$ :  $\phi(s) = \text{nodes}(\text{CART}(s, \text{max\_nodes} = B))$ . The number of MDP states would be  $O((2B)^D)$ . When  $B = 1$ , the MDP policy corresponds to the greedy tree. The process of generating split nodes with calls to CART is illustrated in Figure 6.2.

### 6.3.3 Dynamic Programming to solve the MDP

---

#### Algorithm 4 : DPDT

---

**Data :** Dataset  $\mathcal{E}$ , max depth  $D$ , split function  $\phi()$ ,  
split function parameter  $B$ , regularizing term  $\alpha$   
**Result :** Tree  $T$   
 $\mathcal{M} \leftarrow \text{build\_mdp}(\mathcal{E}, D, \phi(), B)$   
// Backward induction  
 $Q^*(s, a) \leftarrow R_\alpha^\mathcal{M}(s, a) + \sum_{s'} P^\mathcal{M}(s, a, s') \max_{a' \in A_{s'}^\mathcal{M}} Q^*(s', a') \forall s, a \in \mathcal{M}$   
// Get the optimal policy  
 $\pi^*(s) = \arg\max_{a \in A_s^\mathcal{M}} Q^*(s, a) \forall s \in \mathcal{M}$   
// Extracting tree from policy  
 $T \leftarrow E(\pi^*, s_0^\mathcal{M})$

---

After constructing the MDP with a chosen splits generating function, we solve for the optimal policy using dynamic programming. Starting from terminal states and working backward to the initial state, we compute the optimal state-action values using Bellman's

optimality equation [BELLMAN1958228], and then deducing the optimal policy.

From now on, we write DPDT to denote Algorithm 9 when the split function is a call to CART. We discuss key bottlenecks when implementing DPDT in subsequent sections. We now state theoretical results when using DPDT with the CART heuristic.

### 6.3.4 Performance Guarantees of DPDT

We now show that : 1) DPDT minimizes the loss from Eq. 6.1 at least as well as greedy trees and 2) there exists problems for which DPDT has strictly lower loss than greedy trees. As we restrict the action space at a given state  $s$  to a subset of all possible split nodes, DPDT is not guaranteed to find the tree minimizing Eq. 6.1. However, we are still guaranteed to find trees that are better or equivalent to those induced by CART :

**Theorem 1** (MDP solutions are not worse than the greedy tree). *Let  $\pi^*$  be an optimal deterministic policy of the MDP, where the action space at every state is restricted to the top  $B$  most informative or discriminative splits. Let  $T_0$  be the tree induced by CART and  $\{T_1, \dots, T_M\}$  all the sub-trees of  $T_0$ ,<sup>1</sup> then for any  $\alpha > 0$ ,*

$$\mathcal{L}_\alpha(E(\pi^*, s_0)) \leq \min_{0 \leq i \leq M} \mathcal{L}_\alpha(T_i)$$

*Démonstration.* Let us first define  $C(T)$ , the expected number of splits performed by tree  $T$  on dataset  $\mathcal{E}$ . Here  $T$  is deduced from policy  $\pi$ , i.e.  $T = E(\pi, s_0)$ .  $C(T)$  can be defined recursively as  $C(T) = 0$  if  $T$  is a leaf node, and  $C(T) = 1 + p_l C(T_l) + p_r C(T_r)$ , where  $T_l = E(\pi, s_l)$  and  $T_r = E(\pi, s_r)$ . In words, when the root of  $T$  is a split node, the expected number of splits is one plus the expected number of splits of the left and right sub-trees of the root node.  $\square$

It is known that the greedy tree of depth 2 fails to perfectly classify the XOR problem as shown in Figure 5.1 and in [74, 73]. We aim to show that DPDT is a cheap way to alleviate the weaknesses of greedy trees in this type of problems. The following theorem states that there exist classification problems such that DPDT optimizes the regularized training loss strictly better than greedy algorithms such as CART, ID3 or C4.5.

**Theorem 2** (DPDT can be strictly better than greedy). *There exists a dataset and a depth  $D$  such that the DPDT tree  $T_D^{DPDT}$  is strictly better than the greedy tree  $T_D^{greedy}$ , i.e.,  $\mathcal{L}_{\alpha=0}(T_D^{greedy}) > \mathcal{L}_{\alpha=0}(T_D^{DPDT})$ .*

---

1. These sub-trees are interesting to consider since they can be returned by common postprocessing operations following a call to CART, that prune some of the nodes from  $T_0$ . Please see [34] for a review of pruning methods for decision trees.

The proof of this theorem is given in the next section.

### 6.3.5 Proof of Improvement over CART

In this section we construct a dataset for which the greedy tree of depth 2 fails to accurately classify data while DPDT with calls to CART as a splits generating function guarantees a strictly better accuracy. The dataset is the XOR pattern like in Figure 5.1. We will first show that greedy tree induction like CART chooses the first split at random and the second split in between the two columns or rows. Then we will quantify the misclassification of the depth-2 greedy tree on the XOR gate. Finally we will show that using the second greedy split as the root of a tree and then building the remaining nodes greedily, i.e. running DPDT with the CART heuristic, strictly decreases the misclassification.

**Definition 7** (XOR dataset). *Let us defined the XOR dataset as  $\mathcal{E}_{XOR} = \{(X_i, Y_i)\}_{i=1}^N$ .  $X_i = (x_i, y_i) \sim \mathcal{U}([0, 1]^2)$  are i.i.d 2-features samples.  $Y_i = f(X_i)$  are alternating classes with  $f(x, y) = (\lfloor 2x \rfloor + \lfloor 2y \rfloor) \bmod 2$ .*

**Lemma 1.** *The first greedy split is chosen at random on the XOR dataset from definition 13.*

*Démonstration.* Let us consider an arbitrary split  $x = x_v$  parallel to the y-axis. The results apply to splits parallel to the x-axis because the XOR pattern is the same when rotated 90 degrees. The split  $x_v$  partitions the dataset into two regions  $R_{left}$  and  $R_{right}$ . Since the dataset has two columns and two rows, any rectangular area that spans the whole height  $[0, 1)$  has the same proportion of class 0 samples and class 1 samples from definition 13. So in both  $R_{left}$  and  $R_{right}$  the probabilities of observing class 0 or class 1 at random are  $\frac{1}{2}$ . Since the class distributions in left and right regions are independent of the split location, all splits have the same objective value when the objective is a measure of information gain like the entropy or the Gini impurity. Hence, the first split in a greedily induced tree is chosen at random.  $\square$

**Lemma 2.** *When the first split is greedy on the XOR dataset from definition 13, the second greedy splits are chosen perpendicular to the first split at  $y = \frac{1}{2}$*

*Démonstration.* Assume without loss of generality due to symmetries, that the first greedy split is vertical, at  $x = x_v$ , with  $x_v \leq \frac{1}{2}$ . This split partitions the unit square into  $R_{left} = [0, x_v) \times [0, 1)$  and  $R_{right} = [x_v, 1) \times [0, 1)$ . The split  $y = \frac{1}{2}$  further partitions  $R_{left}$  into  $R_{left-down}$  and  $R_{left-up}$  with same areas  $x_v \times y = \frac{x_v}{2}$ . Due to the XOR pattern, there

are only samples of class 0 in  $R_{left-down}$  and only samples of class 1 in  $R_{left-up}$ . Hence the the split  $y = \frac{1}{2}$  maximizes the information gain in  $R_{left}$ , hence the second greedy split given an arbitrary first split  $x = x_v$  is necessarily  $y = \frac{1}{2}$ .  $\square$

**Definition 8** (Forced- Tree). *Let us define the forced-tree as a greedy tree that is forced to make its first split at  $y = \frac{1}{2}$ .*

**Lemma 3.** *The forced-tree of depth 2 has a 0 loss on the XOR dataset from definition 13 while, with probability  $1 - \frac{1}{|\mathcal{E}_{XOR}|}$ , the greedy tree of depth 2 has strictly positive loss.*

*Démonstration.* This is trivial from the definition of the forced tree since if we start with the split  $y = \frac{1}{2}$ , then clearly CART will correctly split the remaining data. If instead the first split is some  $x_v \neq \frac{1}{2}$  then CART is bound to make an error with only one extra split allowed. Since the first split is chosen at random, from Lemma 6.3.5, there are only two splits ( $x = \frac{1}{2}$  and  $y = \frac{1}{2}$ ) out of  $2|\mathcal{E}_{XOR}|$  that do not lead to sub-optimality.  $\square$

We can now formally prove theorem 2.

*Démonstration.* By definition of DPDT, all instances of DPDT with the CART nodes parameter  $B \geq 2$  include the forced-tree from definition 14 in their solution set when applied to the XOR dataset (definition 13). We know from lemma 3 that with high probability, the forced-tree of depth 2 is strictly more accurate than the greedy tree of depth 2 on the XOR dataset. Because we know by proposition 2 that DPDT returns the tree with maximal accuracy from its solution set, we can say that DPDT depth-2 trees are strictly better than depth-2 greedy trees returned by e.g. CART on the XOR dataset.  $\square$

### 6.3.6 Practical Implementation

The key bottlenecks lie in the MDP construction step of DPDT (Section 6.2). In nature, all decision tree induction algorithms have time complexity exponential in the number of training subsets per tree depth  $D : O((2B)^D)$ , e.g., CART has  $O(2^D)$  time complexity. We already saw that DPDT saves time by not considering all possible tree splits but only  $B$  of them. Using state-dependent split generation also allows to generate more or less candidates at different depths of the tree. Indeed, the MDP state  $s = (X, d)$  contains the current depth during the MDP construction process. This means that one can control DPDT's time complexity by giving multiple values of maximum nodes : given  $(B_1, B_2, \dots, B_D)$ , the splits generating function in Algorithm 9

becomes  $\phi(s_i) = \phi(X_i, d = 1) = \text{nodes}(\text{CART}(s, \text{max\_nodes} = B_1))$  and  $\phi(s_j) = \phi(X_j, d = 2) = \text{nodes}(\text{CART}(s, \text{max\_nodes} = B_2))$ .

Similarly, the space complexity of DPDT is exponential in the space required to store training examples  $\mathcal{E}$ . Indeed, the MDP states that DPDT builds in Algorithm 9 are training samples  $X \subseteq \mathcal{E}$ . Hence, the total space required to run DPDT is  $O(Np(2B)^D)$  where  $Np$  is the size of  $\mathcal{E}$ . In practice, one should implement DPDT in a depth first search manner to obtain a space complexity linear in the size of training set :  $O(DNp)$ . In practice DPDT builds the MDP from Section 6.2 by starting from the root and recursively splitting the training set while backpropagating the  $Q$ -values. This is possible because the MDP we solve has a (not necessarily binary) tree structure (see Figure ??) and because the  $Q$ -values of a state only depend on future states.

We implemented DPDT<sup>2</sup> following scikit-learn API [16] with depth-first search and state-depth-dependent splits generating.

## 6.4 Empirical Evaluation

In this section, we empirically demonstrate strong properties of DPDT trees. The first part of our experiments focuses on the quality of solutions obtained by DPDT for objective Eq.6.1 compared to greedy and optimal trees. We know by theorems 1 and 2 that DPDT trees should find better solutions than greedy algorithms for certain problems; but what about real problems? After showing that DPDT can find optimal trees by considering much less solutions and thus performing orders of magnitude less operations, we will study the generalization capabilities of the latter : do DPDT trees label unseen data accurately?

### 6.4.1 DPDT optimizing capabilities

From an empirical perspective, it is key to evaluate DPDT training accuracy since optimal decision tree algorithms against which we wish to compare ourselves are designed to optimize the regularized training loss Eq.6.1.

#### Setup

**Metrics :** we are interested in the regularized training loss of algorithms optimizing Eq.6.1 with  $\alpha = 0$  and a maximum depth  $D$ . We are also interested in the number of key operations performed by each baseline, namely computing candidate split nodes for

2. <https://github.com/KohlerHECTOR/DPDTreeEstimator>

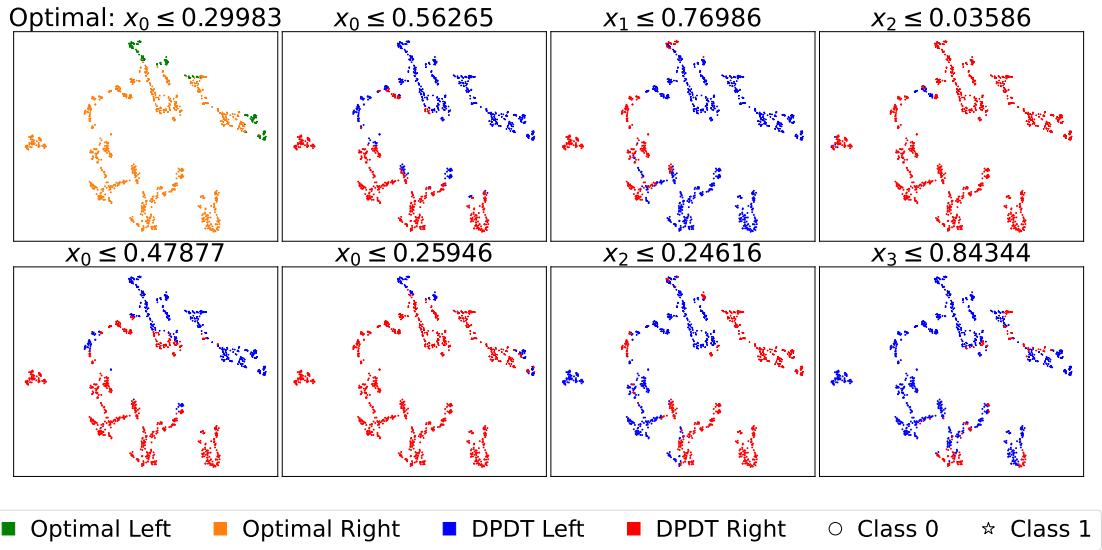


FIGURE 6.3 – Root splits candidate obtained with DPDT compared to the optimal root split on the Bank dataset. Each split creates a partition of  $p$ -dimensional data that we projected in the 2-dimensional space using t-SNE.

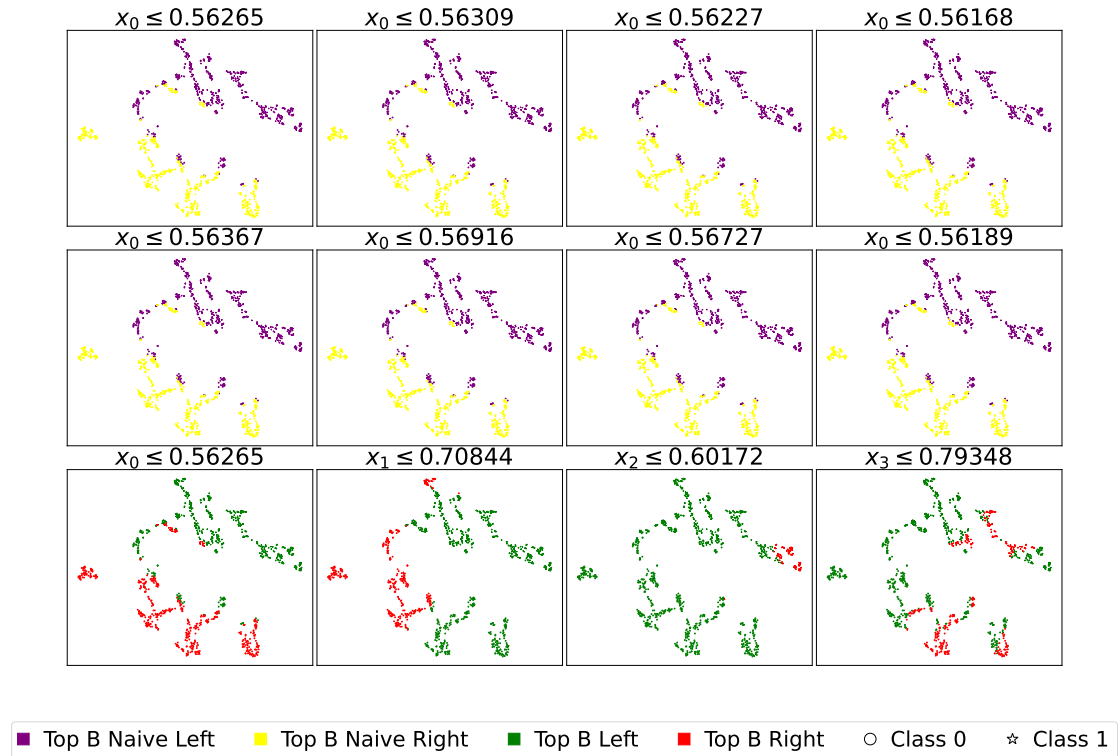


FIGURE 6.4 – Root splits candidate obtained with Top-B[11] on the Bank dataset. Each split creates a partition of  $p$ -dimensional data that we projected using t-SNE.

TABLEAU 6.1 – Comparison of train accuracies of depth-3 trees and number of operations on classification tasks. For DPDT and Top-B, “light” configurations have split function parameters (8, 1, 1) “full” have parameters (8, 8, 8). We also include the mean train accuracy over 5 deep RL runs. **Bold** values are optimal accuracies and blue values are the largest non-optimal accuracies.

Dataset	N	p	Accuracy							Opt Quant-BnB	Greedy CART
			Opt Quant-BnB	Greedy CART	DPDT light	DPDT full	Top-B light	Top-B full	Deep RL Custard		
room	8103	16	<b>0.992</b>	0.968	<b>0.991</b>	<b>0.992</b>	0.990	<b>0.992</b>	0.715	$10^6$	15
bean	10888	16	<b>0.871</b>	0.777	0.812	<b>0.853</b>	0.804	0.841	0.182	$5 \cdot 10^6$	15
eeg	11984	14	<b>0.708</b>	0.666	0.689	<b>0.706</b>	0.684	0.699	0.549	$2 \cdot 10^6$	13
avila	10430	10	<b>0.585</b>	0.532	<b>0.574</b>	<b>0.585</b>	0.563	0.572	0.409	$3 \cdot 10^7$	9
magic	15216	10	<b>0.831</b>	0.801	0.822	<b>0.828</b>	0.807	0.816	0.581	$6 \cdot 10^6$	15
htru	14318	8	<b>0.981</b>	0.979	0.979	<b>0.980</b>	0.979	<b>0.980</b>	0.860	$6 \cdot 10^7$	15
occup.	8143	5	<b>0.994</b>	0.989	0.991	<b>0.994</b>	0.990	<b>0.992</b>	0.647	$7 \cdot 10^5$	13
skin	196045	3	<b>0.969</b>	<b>0.966</b>	<b>0.966</b>	<b>0.966</b>	<b>0.966</b>	<b>0.966</b>	0.612	$7 \cdot 10^4$	15
fault	1552	27	<b>0.682</b>	0.553	0.672	<b>0.674</b>	0.672	0.673	0.303	$9 \cdot 10^8$	13
segment	1848	18	<b>0.887</b>	0.574	0.812	<b>0.879</b>	0.786	0.825	0.137	$2 \cdot 10^6$	7
page	4378	10	<b>0.971</b>	0.964	<b>0.970</b>	<b>0.970</b>	0.964	0.965	0.902	$10^7$	15
bidding	5056	9	<b>0.993</b>	0.981	<b>0.985</b>	<b>0.993</b>	0.985	<b>0.993</b>	0.810	$3 \cdot 10^5$	13
raisin	720	7	<b>0.894</b>	0.869	0.879	<b>0.886</b>	0.875	0.883	0.509	$4 \cdot 10^6$	15
rice	3048	7	<b>0.938</b>	0.933	0.934	<b>0.937</b>	0.933	0.936	0.519	$2 \cdot 10^7$	15
wilt	4339	5	<b>0.996</b>	0.993	0.994	<b>0.995</b>	0.994	0.994	0.984	$3 \cdot 10^5$	13
bank	1097	4	<b>0.983</b>	0.933	0.971	<b>0.980</b>	0.951	0.974	0.496	$6 \cdot 10^4$	13

subsets of the training data. We disregard running times as solvers are implemented in different programming languages and/or using optimized code : operations count is more representative of an algorithm efficiency. We also qualitatively compare different decision trees root splits to some optimal root split.

**Baselines :** we benchmark DPDT against greedy trees and optimal trees. For greedy trees we compare DPDT to CART [13]. For optimal trees we compare DPDT to Quant-BnB [68] which is the only solver specialized for depth 3 trees and continuous features. We also consider the non-greedy baseline Top-B [11]. Ideally, DPDT should have training accuracy close to the optimal tree while performing a number of operations close to the greedy algorithm. Furthermore, comparing DPDT to Top-B brings answers to which heuristic splits are better to consider.

We use the CART algorithm implemented in `scikit-learn` [77] in CPython with a maximum depth of 3. Optimal trees are obtained by running the Julia implementation of the Quant-BnB solver from [68] specialized in depth 3 trees for datasets with

continuous features. We use a time limit of 24 hours per dataset. DPDT and Top-B trees are obtained with Algorithm 9 implemented in pure Python and the calls to CART and Top-B most informative splits generating functions from Section 6.2 respectively. We also include Custard, a deep RL baseline [97]. Custard fits a neural network online one datum at a time rather than solving exactly the MDP from Section 6.2 which states are sets of data. Similarly to DPDT, Custard neural network policy is equivalent to a decision tree. We implement Custard with the DQN agent from `stable-baselines3` [85] and train until convergence.

**Datasets :** we use the same datasets as the Quant-BnB paper [68].

### Observations

**Near-optimality.** Our experimental results demonstrate that unlike Deep RL, DPDT and Top-B approaches consistently improve upon greedy solutions while requiring significantly fewer operations than exact solvers. Looking at Table 6.1, we observe several key patterns : first, light DPDT with 16 candidate root splits consistently outperforms the greedy baseline in all datasets. This shows that in practice DPDT can be strictly better than CART outside of theorem 2 assumptions. Second, when comparing DPDT to Top-B, we see that DPDT generally achieves better accuracy for the same configuration. For example, on the bean dataset, full DPDT reaches 85.3% accuracy while full Top-B achieves 84.1%. This pattern holds on most datasets, suggesting that DPDT is more effective than selecting splits based purely on information gain.

Third, both approaches achieve impressive computational efficiency compared to exact solvers. While optimal solutions require between  $10^4$  to  $10^8$  operations, DPDT and Top-B typically need only  $10^2$  to  $10^4$  operations, a reduction of 2 to 4 orders of magnitude. Notably, on several datasets (room, avila, occupancy, bidding), full DPDT matches or comes extremely close to optimal accuracy while requiring far fewer operations. For example, on the room dataset, full DPDT achieves the optimal accuracy of 99.2% while reducing operations from  $1.34 \times 10^6$  to  $1.61 \times 10^4$ . These results demonstrate that DPDT provides an effective middle ground between greedy approaches and exact solvers, offering near-optimal solutions with reasonable computational requirements. While both DPDT and Top-B improve upon greedy solutions, DPDT CART-based split generation strategy appears to be particularly effective at finding high-quality solutions.

**DPDT splits** To understand why the CART-based split generation yields more accurate DPDT trees than the Top-B heuristic, we visualize how splits partition the feature space



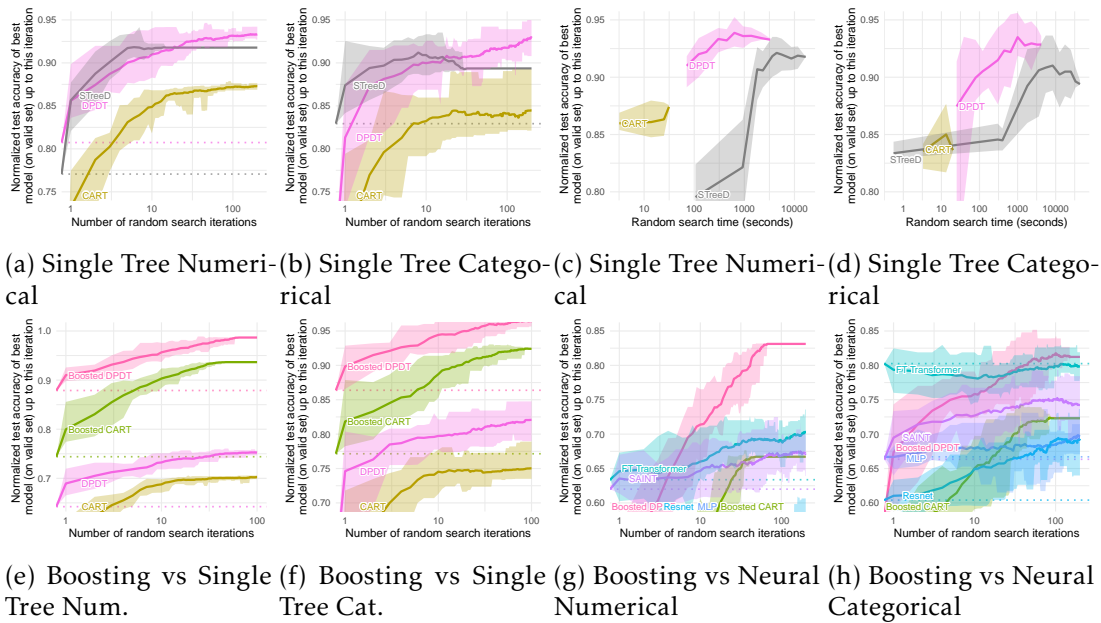


FIGURE 6.5 – Benchmark on medium-sized datasets. Dotted lines correspond to the score of the default hyperparameters, which is also the first random search iteration. Each value corresponds to the test score of the best model (obtained on the validation set) after a specific number of random search iterations (a, b) or after a specific time spent doing random search (c, d), averaged on 15 shuffles of the random search order. The ribbon corresponds to the minimum and maximum scores on these 15 shuffles.

(Figures 6.3, 6.4). We run both DPDT with splits from CART and DPDT with the Top-B most informative splits on the bank dataset. We use t-SNE to create a two-dimensional representations of the dataset partitions given by candidates root splits from CART and Top-B. The optimal root split for the depth-3 tree for bank—obtained with Quant-BnB—is shown on Figure 6.3 in the top-left subplot using green and orange colors for the resulting partitions. On the same figure we can see that the DPDT split generated with CART  $x_0 \leq 0.259$  is very similar to the optimal root split. However, on Figure 6.4 we observe that no Top-B candidate splits resemble the optimal root and that in general Top-B split lack diversity : they always split along the same feature. We tried to enforce diversity by getting the most informative split *per feature* but no candidate split resembles the optimal root.

### 6.4.2 DPDT generalization capabilities

The goal of this section is to have a fair comparison of generalization capabilities of different tree induction algorithms. Fairness of comparison should take into account the number of hyperparameters, choice of programming language, intrinsic purposes of each algorithms (what are they designed to do?), the type of data they can read (categorical features or numerical features). We benchmark DPDT using [44]. We choose this benchmark because it was used to establish XGBoost [20] as the SOTA tabular learning model.

#### Setup

**Metrics :** We re-use the code from [44]<sup>3</sup>. It relies on random searches for hyperparameter tuning [8]. We run a random search of 100 iterations per dataset for each benchmarked tree algorithms. To study performance as a function of the number  $n$  of random search iterations, we compute the best hyperparameter combination on the validation set on these  $n$  iterations (for each model and dataset), and evaluate it on the test set. Following [44], we do this 15 times while shuffling the random search order at each time. This gives us bootstrap-like estimates of the expected test score of the best tree found on the validation set after each number of random search iterations. In addition, we always start the random searches with the default hyperparameters of each tree induction algorithm. We use the test set accuracy (classification) to measure model performance. The aggregation metric is discussed in details in [44, Section 3].

**Datasets :** we use the datasets curated by [44]. They are available on OpenML [101] and described in details in [44, Appendix A.1]. The attributes in these datasets are either numerical (a real number), or categorical (a symbolic values among a finite set of possible values). The considered datasets follow a strict selection [44, Section 3] to focus on core learning challenges. Some datasets are very large (millions of samples) like Higgs or Coverttype [109, 10]. To ensure non-trivial learning tasks, datasets where simple models (e.g. logistic regression) performed within 5% of complex models (e.g. ResNet [43], HistGradientBoosting [77]) are removed. We use the same data partitioning strategy as [44] : 70% of samples are allocated for training, with the remaining 30% split between validation (30%) and test (70%) sets. Both validation and test sets are capped at 50,000 samples for computational efficiency. All algorithms and hyperparameter combinations were evaluated on identical folds. Finally, while we focus on classification

3. <https://github.com/leogrin/tabular-benchmark>

TABLEAU 6.2 – Hyperparameters importance comparison. A description of the hyperparameters can be found in the scikit-learn documentation : <https://scikit-learn.org/stable/modules/generated/sklearn.tree.DecisionTreeClassifier.html>.

Hyperparameter	DPDT (%)	CART (%)	STreeD (%)
min_samples_leaf	35.05	33.50	50.50
min_impurity_decrease	24.60	24.52	-
cart_nodes_list	15.96	-	-
max_features	11.16	18.06	-
max_depth	7.98	10.19	0.00
max_leaf_nodes	-	7.84	-
min_samples_split	2.67	2.75	-
min_weight_fraction_leaf	2.58	3.14	-
max_num_nodes	-	-	27.51
n_thresholds	-	-	21.98

datasets in the main text, we provide results for regression problems in table 6.5 in the appendix.

**Baselines :** we benchmark DPDT against CART and STreeD when inducing trees of depth at most 5. We use hyperparameter search spaces from [55] for CART and DPDT. For DPDT we additionally consider eight different splits functions parameters configurations for the maximum nodes in the calls to CART. Surprisingly, after computing the importance of each hyperparameter of DPDT, we found that the maximum node numbers in the calls to CART are only the third most important hyperparameter behind classical ones like the minimum size of leaf nodes or the minimum impurity decrease (Table 6.2). We use the CPython implementation of STreeD<sup>4</sup>. All hyperparameter grids are given in table 6.7 in the appendix.

**Hardware :** experiments were conducted on a heterogeneous computing infrastructure made of AMD EPYC 7742/7702 64-Core and Intel Xeon processors, with hardware allocation based on availability and algorithm requirements. DPDT and CART random searches ran for the equivalent of 2-days while PySTreeD ran for 10-days.

## Observations

**Generalization** In Figure 6.5, we observe that DPDT learns better trees than CART and STreeD both in terms of generalization capabilities and in terms of computation

4. PySTreeD : <https://github.com/AlgTUDelft/pystreed>

TABLEAU 6.3 – Depth-10 decision trees for the KDD 1999 cup dataset.

Model	Test Accuracy (%)	Time (s)	Memory (MB)
DPDT-(4,)	<b>91.30</b>	339.85	5054
DPDT-(4,4,)	<b>91.30</b>	881.07	5054
CART	91.29	25.36	1835
GOSDT- $\alpha = 0.0005$	65.47	5665.47	1167
GOSDT- $\alpha = 0.001$	65.45	5642.85	1167

cost. On Figures 6.5a and 6.5b, DPDT obtains best generalization scores for classification on numerical and categorical data after 100 iterations of random hyperparameters search over both CART and STreeD. Similarly, we also present generalization scores as a function of compute time (instead of random search iterations). On Figures 6.5c and 6.5d, despite being coded in the slowest language (Python vs. CPython), our implementation of DPDT finds the best overall model before all STreeD random searches even finish. The results from Figure 6.5 are appealing for machine learning practitioners and data scientists that have to do hyperparameters search to find good models for their data while having computation constraints.

Now that we have shown that DPDT is extremely efficient to learn shallow decision trees that generalize well to unseen data, it is fair to ask if DPDT can also learn deep trees on very large datasets.

**Deeper trees on bigger datasets.** We also stress test DPDT by inducing deep trees of depth 10 for the KDD 1999 cup dataset <sup>5</sup>. The training set has 5 million rows and a mix of 80 continuous and categorical features representing network intrusions. We fit DPDT with 4 split candidates for the root node (DPDT-(4,)) and with 4 split candidates for the root and for each of the internal nodes at depth 1 (DPDT-(4,4,)). We compare DPDT to CART with a maximum depth of 10 and to GOSDT <sup>6</sup> [McTavish\_Zhong\_Achermann\_Karimalis\_Chen] with different regularization values  $\alpha$ . GOSDT first trains a tree ensemble to binarize a dataset and then solve for the optimal decision tree of depth 10 on the binarized problem. In Table 6.3 we report the test accuracy of each tree on the KDD 1999 cup test set. We also report the memory peak during training and the training duration (all experiments are run on the same CPU). We observe that DPDT can improve over CART even for deep trees and large datasets while using reasonable time and memory. Furthermore, Table 6.3 highlights the limitation of optimal trees for practical problems when the dataset is not binary. We observed that GOSDT could not find a good binarization of the dataset

5. <http://kdd.ics.uci.edu/databases/kddcup99/kddcup99.html>

6. Code available at: <https://github.com/ubc-systopia/gosdt-guesses>

even when increasing the budget of the tree ensemble up to the point where most of the computations are spent on fitting the ensemble (see more details about this phenomenon in [McTavish\_Zhong\_Achermann\_Karimalis\_Chen\_Rudin\_Seltzer\_2022]). In table 6.6 in the appendix, we also show that DPDT performs better than optimal trees for natively binary datasets. In the next section we study the performance of boosted DPDT trees.

## 6.5 Application of DPDT to Boosting

In the race for machine learning algorithms for tabular data, boosting procedures are often considered the go-to methods for classification and regression problems. Boosting algorithms [37, 39, 38] sequentially add weak learners to an ensemble called strong learner. The development of those boosting algorithms has focused on what data to train newly added weak learners [39, 38], or on efficient implementation of those algorithms [20, 79]. We show next that Boosted-DPDT (boosting DPDT trees with AdaBoost [37]) improves over recent deep learning algorithms.

### 6.5.1 Boosted-DPDT

We benchmark Boosted-DPDT with the same datasets, metrics, and hardware as in the previous section on single-tree training. Second, we verify the competitiveness of Boosted-DPDT with other models such as deep learning ones (SAINT [94] and other deep learning architectures from [43]).

On Figures 6.5e and 6.5f we can notice 2 properties of DPDT. First, as in any boosting procedure, Boosted-DPDT outperforms its weak counterpart DPDT. This serves as a sanity check for boosting DPDT trees. Second, it is clear that boosting DPDT trees yields better models than boosting CART trees on both numerical and categorical data. Figures 6.5g and 6.5h show that boosting DPDT trees using the default AdaBoost procedure [37] is enough to learn models outperforming deep learning algorithms on datasets with numerical features and models in the top-tier on datasets with categorical features. This show great promise for models obtained when boosting DPDT trees with more advanced procedures.

### 6.5.2 (X)GB-DPDT

We also boost DPDT trees with Gradient Boosting and eXtreme Gradient Boosting [38, 39, 20](X(GB)-DPDT). For each dataset from [44], we trained (X)GB-DPDT

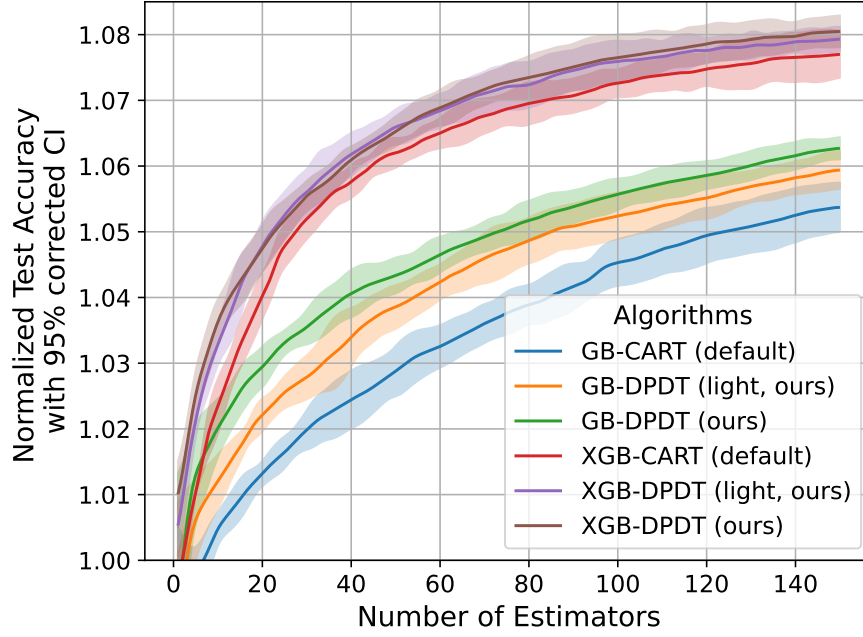


FIGURE 6.6 – Aggregated mean test accuracies of Gradient Boosting models as a function of the number of single trees.

models with 150 boosted single DPDT trees and a maximum depth of 3 for each. We evaluate two DPDT configurations for the single trees : light (DPDT-(4, 1, 1)) and the default (DPDT-(4,4,4)). We compare (X)GB-DPDT to (X)GB-CART : 150 boosted CART trees with maximum depth of 3 and default hyperparameters for each. All models use a learning rate of 0.1. For each dataset, we normalize all boosted models scores by the accuracy of a single depth-3 CART decision tree and aggregate the results : the final curves represent the mean performance across all datasets, with confidence intervals computed using 5 different random seeds.

Figure 6.6 shows that similarly to simple boosting procedures like AdaBoost, more advanced ones like (eXtreme) Gradient Boosting yields better models when the weak learners are DPDT trees rather than greedy trees. This is a motivation to develop efficient implementation of (eXtreme) Gradient Boosting with DPDT as the weak learning algorithm to perform extensive benchmarking following [44] and potentially claim the state-of-the-art.

## 6.6 Proof of Propostion 2

For the purpose of the proof, we overload the definition of  $J_\alpha$  and  $\mathcal{L}_\alpha$ , to make explicit the dependency on the dataset and the maximum depth. As such,  $J_\alpha(\pi)$  becomes  $J_\alpha(\pi, \mathcal{E}, D)$  and  $\mathcal{L}_\alpha(T)$  becomes  $\mathcal{L}_\alpha(T, \mathcal{E})$ . Let us first show that the relation  $J_\alpha(\pi, \mathcal{E}, 0) = -\mathcal{L}_\alpha(T, \mathcal{E})$  is true. If the maximum depth is  $D = 0$  then  $\pi(s_0)$  is necessarily a class assignment, in which case the expected number of splits is zero and the relation is obviously true since the reward is the opposite of the average classification loss. Now assume it is true for any dataset and tree of depth at most  $D$  with  $D \geq 0$  and let us prove that it holds for all trees of depth  $D + 1$ . For a tree  $T$  of depth  $D + 1$  the root is necessarily a split node. Let  $T_l = E(\pi, s_l)$  and  $T_r = E(\pi, s_r)$  be the left and right sub-trees of the root node of  $T$ . Since both sub-trees are of depth at most  $D$ , the relation holds and we have  $J_\alpha(\pi, X_l, D) = \mathcal{L}_\alpha(T_l, X_l)$  and  $J_\alpha(\pi, X_r, D) = \mathcal{L}_\alpha(T_r, X_r)$ , where  $X_l$  and  $X_r$  are the datasets of the “right” and “left” states to which the MDP transitions—with probabilities  $p_l$  and  $p_r$ —upon application of  $\pi(s_0)$  in  $s_0$ , as described in the MDP formulation. Moreover, from the definition of the policy return we have

$$\begin{aligned}
 J_\alpha(\pi, \mathcal{E}, D + 1) &= -\alpha + p_l * J_\alpha(\pi, X_l, D) + p_r * J_\alpha(\pi, X_r, D) \\
 &= -\alpha - p_l * \mathcal{L}_\alpha(T_l, X_l) - p_r * \mathcal{L}_\alpha(T_r, D) \\
 &= -\alpha - p_l * \left( \frac{1}{|X_l|} \sum_{(x_i, y_i) \in X_l} \ell(y_i, T_l(x_i)) + \alpha C(T_l) \right) \\
 &\quad - p_r * \left( \frac{1}{|X_r|} \sum_{(x_i, y_i) \in X_r} \ell(y_i, T_r(x_i)) + \alpha C(T_r) \right) \\
 &= -\frac{1}{N} \sum_{(x_i, y_i) \in X} \ell(y_i, T(x_i)) - \alpha(1 + p_l C(T_l) + p_r C(T_r)) \\
 &= -\mathcal{L}(T, \mathcal{E})
 \end{aligned}$$

## 6.7 Additional Experiments and Hyperparameters

In this section we provide additional experimental results. In Table 6.5, we compare DPDT trees to CART and STreeD trees using 50 train/test splits of regression datasets from [44]. All algorithms are run with default hyperparameters.

The configuration of DPDT is (4, 4, 4) or (4,4,4,4,4). STreeD is run with a time limit of 4 hours per tree computation and on binarized versions of the datasets. Both for depth-3 and depth-5 trees, DPDT outperforms other baselines in terms of train and test accuracies. Indeed, because STreeD runs on “approximated” datasets, it performs pretty poorly.

In Table 6.6, we compare DPDT(4, 4, 4, 4, 4) to additional optimal decision tree baselines on datasets with **binary features**. The optimal decision tree baselines run with default hyperparameters and a time-limit of 10 minutes. The results show that even on binary datasets that optimal algorithms are designed to handle well; DPDT outperforms other baselines. This is likely because optimal trees are slow and/or don’t scale well to depth 5.

In Table 6.4 compare DPDT to lookahead depth-3 trees when optimizing Eq.6.1. Unlike the other greedy approaches, lookahead decision trees [**norton**] do not pick the split that optimizes a heuristic immediately. Instead, they pick a split that sets up the best possible heuristic value on the following split. Lookahead-1 chooses nodes at depth  $d < 3$  by looking 1 depth in the future : it looks for the sequence of 2 splits that maximizes the information gain at depth  $d + 1$ . Lookahead-2 is the optimal depth-3 tree and Lookahead-0 would be just building the tree greedily like CART. The conclusion are roughly the same as for Table 6.1. Both lookahead trees and DPDT<sup>7</sup> are in Python which makes them slow but comparable.

We also provide the hyperparameters necessary to reproduce experiments from section 6.4.2 and 6.5.1 in Table 6.7.

---

7. <https://github.com/KohlerHECTOR/DPDTreeEstimator>



TABLEAU 6.4 – Train accuracies of depth-3 trees (with number of operations). Lookahead trees are trained with a time limit of 12 hours.

Dataset	DPDT	Lookahead-1
avila	57.22 (1304)	OoT
bank	97.99 (699)	96.54 (7514)
bean	85.30 (1297)	OoT
bidding	99.27 (744)	98.12 (20303)
eeg	69.38 (1316)	69.09 (10108)
fault	67.40 (1263)	67.20 (32514)
htru	98.01 (1388)	OoT
magic	82.81 (1451)	OoT
occupancy	99.31 (1123)	99.01 (15998)
page	97.03 (1243)	96.44 (16295)
raisin	88.61 (1193)	86.94 (9843)
rice	93.44 (1367)	93.24 (37766)
room	99.23 (1196)	99.04 (5638)
segment	87.88 (871)	68.83 (24833)
skin	96.60 (1300)	96.61 (1290)
wilt	99.47 (862)	99.31 (36789)

TABLEAU 6.5 – Mean train and test scores (with standard errors) for regression datasets over 50 cross-validation runs.

	Depth 3								
Dataset	DPDT		Optimal		CART		DPDT		
	Train	Test	Train	Test	Train	Test	Train	Test	
nyc-taxi	39.0 ± 0.0	38.9 ± 0.2	33.8 ± 0.0	33.8 ± 0.1	39.0 ± 0.0	38.9 ± 0.2	45.8 ± 0.0	45.7 ± 0.2	33.
medical_charges	95.2 ± 0.0	95.2 ± 0.0	90.1 ± 0.0	90.1 ± 0.1	95.2 ± 0.0	95.2 ± 0.0	97.7 ± 0.0	97.7 ± 0.0	90.
diamonds	93.0 ± 0.0	92.9 ± 0.1	90.1 ± 0.0	90.1 ± 0.1	92.7 ± 0.0	92.6 ± 0.1	94.2 ± 0.0	94.0 ± 0.1	90.
house_16H	39.9 ± 0.1	38.1 ± 2.5	32.8 ± 0.1	29.4 ± 1.6	35.8 ± 0.1	35.8 ± 1.9	59.4 ± 0.1	35.2 ± 4.1	32.
house_sales	67.0 ± 0.0	66.0 ± 0.4	65.1 ± 0.0	64.4 ± 0.4	66.8 ± 0.0	66.1 ± 0.4	77.6 ± 0.0	76.1 ± 0.3	65.
superconduct	73.1 ± 0.0	72.7 ± 0.5	70.9 ± 0.0	70.5 ± 0.5	70.4 ± 0.0	69.7 ± 0.5	83.0 ± 0.0	81.7 ± 0.4	70.
houses	51.7 ± 0.0	50.7 ± 0.7	48.5 ± 0.1	47.3 ± 0.7	49.5 ± 0.0	48.4 ± 0.7	69.1 ± 0.0	67.6 ± 0.5	48.
Bike_Sharing	55.2 ± 0.0	54.7 ± 0.5	45.1 ± 0.1	44.8 ± 0.7	48.1 ± 0.0	47.9 ± 0.5	65.2 ± 0.0	63.3 ± 0.5	45.
elevators	48.0 ± 0.0	46.8 ± 1.1	40.2 ± 0.1	38.2 ± 1.0	46.8 ± 0.0	45.5 ± 1.2	65.6 ± 0.0	61.2 ± 1.0	40.
pol	72.2 ± 0.0	71.3 ± 0.6	67.8 ± 0.1	67.5 ± 0.9	72.0 ± 0.0	71.2 ± 0.8	93.3 ± 0.0	92.4 ± 0.3	67.
MiamiHousing2016	62.3 ± 0.0	60.4 ± 0.8	60.8 ± 0.0	58.3 ± 0.8	62.3 ± 0.0	60.6 ± 0.8	79.8 ± 0.0	77.5 ± 0.5	60.
Ailerons	63.5 ± 0.0	62.6 ± 0.7	61.6 ± 0.0	60.3 ± 0.7	63.5 ± 0.0	62.6 ± 0.7	76.0 ± 0.0	72.9 ± 0.6	61.
Brazilian_houses	90.7 ± 0.0	90.3 ± 0.8	89.2 ± 0.0	89.4 ± 0.8	90.7 ± 0.0	90.4 ± 0.8	97.6 ± 0.0	96.6 ± 0.4	89.
sulfur	72.5 ± 0.1	66.6 ± 2.2	35.7 ± 0.1	19.1 ± 6.7	72.0 ± 0.1	68.0 ± 2.2	89.0 ± 0.0	68.4 ± 6.7	35.
yprop_41	6.3 ± 0.0	2.3 ± 0.7	3.6 ± 0.0	1.5 ± 0.4	6.2 ± 0.0	2.1 ± 0.8	13.2 ± 0.0	1.2 ± 1.7	3.
cpu_act	93.4 ± 0.0	92.0 ± 0.6	89.0 ± 0.0	86.5 ± 1.9	93.4 ± 0.0	92.0 ± 0.6	96.5 ± 0.0	94.7 ± 0.5	89.
wine_quality	27.9 ± 0.0	23.3 ± 0.9	25.2 ± 0.0	23.7 ± 0.8	27.7 ± 0.0	24.5 ± 0.8	37.4 ± 0.0	26.7 ± 1.0	25.
abalone	46.3 ± 0.0	39.6 ± 1.6	42.5 ± 0.0	40.4 ± 1.4	43.3 ± 0.0	39.2 ± 1.2	58.6 ± 0.0	44.7 ± 1.8	42.

TABLEAU 6.6 – Train/test accuracies of different decision tree induction algorithms. All algorithms induce trees of depth at most 5 on 8 classification datasets. A time limit of 10 minutes is set for OCT-type algorithms. The values in this table are averaged over 3 seeds giving 3 different train/test datasets.

Names	Datasets			Train Accuracy depth-5					Test	
	Samples	Features	Classes	DPDT	OCT	MFOCT	BinOCT	CART	DPDT	OCT
balance-scale	624	4	3	90.9%	71.8%	82.6%	67.5%	86.5%	77.1%	66.9%
breast-cancer	276	9	2	94.2%	88.6%	91.1%	75.4%	87.9%	66.4%	67.1%
car-evaluation	1728	6	4	92.2%	70.1%	80.4%	84.0%	87.1%	90.3%	69.5%
hayes-roth	160	9	3	93.3%	82.9%	95.4%	64.6%	76.7%	75.4%	77.5%
house-votes-84	232	16	2	100.0%	100.0%	100.0%	100.0%	99.4%	95.4%	93.7%
soybean-small	46	50	4	100.0%	100.0%	100.0%	76.8%	100.0%	93.1%	94.4%
spect	266	22	2	93.0%	92.5%	93.0%	92.2%	88.5%	73.1%	75.6%
tic-tac-toe	958	24	2	90.8%	68.5%	76.1%	85.7%	85.8%	82.1%	69.6%

TABLEAU 6.7 – Hyperparameter search spaces for tree-based models. More details about the hyperparameters meaning are given in [55].

Parameter	CART	Boosted-CART	DPDT	Boosted-DPDT	STreeD
<i>Common Tree Parameters</i>					
max_depth	{5 : 0.7, 2,3,4 : 0.1}	{2 : 0.4, 3 : 0.6}	{5 : 0.7, 2,3,4 : 0.1}	{2 : 0.4, 3 : 0.6}	5
min_samples_split	{2 : 0.95, 3 : 0.05}	{2 : 0.95, 3 : 0.05}	{2 : 0.95, 3 : 0.05}	{2 : 0.95, 3 : 0.05}	–
min_impurity_decrease	{0.0 : 0.85, 0.01,0.02,0.05 : 0.05}	{0.0 : 0.85, 0.01,0.02,0.05 : 0.05}	{0.0 : 0.85, 0.01,0.02,0.05 : 0.05}	{0.0 : 0.85, 0.01,0.02,0.05 : 0.05}	–
min_samples_leaf	$\mathcal{Q}(\log -\mathcal{U}[2,51])$	$\mathcal{Q}(\log -\mathcal{U}[2,51])$	$\mathcal{Q}(\log -\mathcal{U}[2,51])$	$\mathcal{Q}(\log -\mathcal{U}[2,51])$	$\mathcal{Q}(\log -\mathcal{U}[2,51])$
min_weight_fraction_leaf	{0.0 : 0.95, 0.01 : 0.05}	{0.0 : 0.95, 0.01 : 0.05}	{0.0 : 0.95, 0.01 : 0.05}	{0.0 : 0.95, 0.01 : 0.05}	–
max_features	{ "sqrt" : 0.5, "log2" : 0.25, 10000 : 0.25 }	{ "sqrt" : 0.5, "log2" : 0.25, 10000 : 0.25 }	{ "sqrt" : 0.5, "log2" : 0.25, 10000 : 0.25 }	{ "sqrt" : 0.5, "log2" : 0.25, 10000 : 0.25 }	–
<i>Model-Specific Parameters</i>					
max_leaf_nodes	{32 : 0.85, 5,10,15 : 0.05}	{8 : 0.85, 5 : 0.05, 7 : 0.1}	–	–	–
cart_nodes_list	–	–	8 configs (uniform)	5 configs (uniform)	–
learning_rate	–	$\log \mathcal{N}(\ln(0.01), \ln(10))$	–	$\log \mathcal{N}(\ln(0.01), \ln(10))$	–
n_estimators	–	1000	–	1000	–
max_num_nodes	–	–	–	–	{3,5,7,11, 17,25,31} (uniform)
n_thresholds	–	–	–	–	{5,10,20,50} (uniform)
cost_complexity	–	–	–	–	0
time_limit	–	–	–	–	1800

# Conclusion

## 7.1 Conclusion

In this paper, we introduced Dynamic Programming Decision Trees (DPDT), a novel framework that bridges the gap between greedy and optimal decision tree algorithms. By formulating tree induction as an MDP and employing adaptive split generation based on CART, DPDT achieves near-optimal training loss with significantly reduced computational complexity compared to existing optimal tree solvers. Furthermore, we prove that DPDT can learn strictly more accurate trees than CART.

Most importantly, extensive benchmarking on varied large and difficult enough datasets showed that DPDT trees and boosted DPDT trees generalize better than other baselines. To conclude, we showed that DPDT is a promising machine learning algorithm.

The key future work would be to make DPDT industry-ready by implementing it in C and or making it compatible with the most advanced implementation of e.g. XGBoost.

## 7.2 What about imitation?



**Troisième partie**

**Beyond Decision Trees : what can be  
done with other Interpretable  
Policies?**



# Imitation and Evaluation

## 8.1 Intro

There exist applications of reinforcement learning like medicine where policies need to be “interpretable” by humans. User studies have shown that some policy classes might be more interpretable than others. However, it is costly to conduct human studies of policy interpretability. Furthermore, there is no clear definition of policy interpretability, i.e., no clear metrics for interpretability and thus claims depend on the chosen definition. We tackle the problem of empirically evaluating policies interpretability without humans. Despite this lack of clear definition, researchers agree on the notions of “*simulatability*” : policy interpretability should relate to how humans understand policy actions given states. To advance research in interpretable reinforcement learning, we contribute a new methodology to evaluate policy interpretability. We distillate expert neural networks policies into small programs that we use as baselines. We then show that using our methodology to evaluate the baselines interpretability leads to similar conclusions as user studies. Most importantly, we show that there is no policy class that better trades off interpretability and performance across tasks.

There is increasing research in developing reinforcement learning algorithms that return “interpretable” policies such as trees, programs, first-order logic, or linear maps [kohler2024interpretableeditableprogrammattictree, 6, 103, 64, 25, 70, 42]. Indeed, interpretability has been useful for different applications : policy verification [6], misalignment detection [26, 67] and features importance analysis [108, 4, 1].

User studies have established the common beliefs that decision trees are more “interpretable” than linear maps, oblique trees (trees where nodes are tests of linear combina-

tions of features), and multi-layer perceptrons (MLPs) [35, 36, 66, 102]. Furthermore, for a fixed class of models, humans give different values of interpretability to models with different numbers of parameters [56]. However, survey works argue that every belief about interpretability needs to be verified with user studies and that interpretability evaluations are grounded to a specific set of users, to a specific application, and to a specific definition of interpretability [**rigorous**, **mythos**]. For example, [**mythos**] claims that depending on the notion of *simulatability* studied, MLPs can be more interpretable than trees, since deep trees can be harder for a human to read than compact MLPs. Hence, even with access to users it would be difficult to research interpretability. More realistically, since the cost of user studies is high (time, variety of subjects required, ethics, etc.), designing proxies for interpretability in machine learning has become an important open problem in both supervised [**rigorous**] and reinforcement learning [42].

In this work, we propose a methodology to evaluate the interpretability of reinforcement learning without human evaluators, by measuring inference times and memory consumptions of policies as programs. We show that those measures constitute adequate proxies for the notions of “*simulatability*” described in [**mythos**], which relates the interpretability of policy to humans ability to understand the inference of actions given states. In addition to the contributions summarized next, we open source some of the interpretable baselines to be used for future interpretability research and teaching<sup>1</sup>.

---

1. <https://anonymous.4open.science/r/interpretable-rl-zoo-4DCC/README.md>



# Evaluation

## 9.1 Methodology Overview

In this section, we explain our methodology for the evaluation of interpretability. Our approach consists of three main steps : (1) obtaining deep neural network policies trained with reinforcement learning that obtain high cumulative rewards, (2) distilling those policies into less complex ones to use as baselines (3) after parsing baselines from different classes into a common comparable language, we evaluate the interpretability of the policies using proxy metrics for *simulatability*.

**Deep Neural Network Policies** In reinforcement learning, an agent learns how to behave in an environment to maximize a cumulative reward signal [95]. The environment is defined by a Markov decision process (MDP)  $M = (\mathcal{S}, \mathcal{A}, T, R)$ , where  $\mathcal{S}$  is the state space,  $\mathcal{A}$  is the action space,  $T : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$  is the state-transition function,  $R : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}$  is the reward function. At each time step  $t$ , an agent observes the current state  $s_t$  and chooses an action  $a_t$  according to its policy  $\pi$ . The agent executes  $a_t$ , receives reward  $r_t$ , and observes the next state  $s'_{t+1}$ . The goal is to find an optimal policy  $\pi^*$  that maximizes the expected discounted future return :  $\pi^* = \operatorname{argmax}_{\pi} Q^{\pi}(s, a) = \operatorname{argmax} \mathbb{E}[r + \gamma Q^{\pi}(s', a)]$ , with  $\gamma$  a discount factor in  $[0, 1]$ . For large or continuous state spaces like the MDPs we consider in this work, MLPs are used to represent  $Q^{\pi}$  or  $\pi$ . While these MLPs can be trained efficiently to obtain high cumulative rewards [90, 71], they are too complex for interpretability considerations.

**Distilling into Interpretable Policies** To obtain interpretable policies, we distill the complex neural networks into simpler models using imitation learning, as described in

**Algorithm 5 : Imitate Expert** [78, 88, 6]

---

**Input :** Expert policy  $\pi^*$ , MDP  $M$ , policy class  $\Pi$ , number of iterations  $N$ , total samples  $S$ , importance sampling flag  $I$

**Output :** Fitted student policy  $\hat{\pi}_i$

Initialize dataset  $\mathcal{D} \leftarrow \emptyset$ ;

Initialize  $\hat{\pi}_1$  arbitrarily from  $\Pi$ ;

**for**  $i \leftarrow 1$  **to**  $N$  **do**

**if**  $i = 1$  **then**  $\pi_i \leftarrow \pi^*$  ;

**else**  $\pi_i \leftarrow \hat{\pi}_i$  ;

Sample  $S/N$  transitions from  $M$  using  $\pi_i$ ;

**if**  $I$  is *True* **then**  $w(s) \leftarrow V^{\pi^*}(s) - \min_a Q^{\pi^*}(s, a)$  ;

**else**  $w(s) \leftarrow 1$  ;

Collect dataset  $\mathcal{D}_i \leftarrow \{(s, \pi^*(s), w(s))\}$  of states visited by  $\pi_i$  and expert actions;

$\mathcal{D} \leftarrow \mathcal{D} \cup \mathcal{D}_i$ ;

Fit classifier/regressor  $\hat{\pi}_{i+1}$  on  $\mathcal{D}$ ;

**end**

**return**  $\hat{\pi}_N$ ;

---

Algorithm 10. This approach transforms the reinforcement learning task into a sequence of supervised learning problems.

Algorithm 10 inputs an environment, that simulates taking steps in an MDP, an expert policy to imitate, also called a teacher, and an (interpretable) policy class to fit, also called student. The hyperparameters of Algorithm 10 are : the number of times we fit a student policy, the total number of samples to be collected, and whether or not to use importance sampling. At each iteration of Algorithm 10 the student policy is fitted to a dataset of states collected with the expert at iteration 1 or with the previously fitted student (see Line 10). The actions are always given by the expert (see Line 10). When using importance sampling, the states are further re-weighted by the worst state-action value possible in the given state. When the number of iteration is 1, Algorithm 10 is behavior cloning [78]. When we use importance sampling, Algorithm 10 is Q-Dagger [6]. In other cases, Algorithm 10 is DAgger [88].

**Measuring Policy Interpretability** After obtaining interpretable policy baselines using Algorithm 10, we use two metrics to evaluate policy interpretability without requiring human intervention. Those metrics are proxies for the notion of *simulatability* from [mythos] that gives insights on how a human being would read a policy to understand how actions are inferred. In particular, *simulatability* admits two sub-definitions. The first one is a measure of how difficult it is for a human to reproduce the computa-

tions of the policy to infer actions given states. The second one measures how difficult it is for a human to read through the entire policy. [mythos] argues that this nuance is key when measuring interpretability because a tree is not read entirely to compute a single action and because there is no consensus on what is easier for a human to read between an MLP and a tree.

1. *Policy Inference Time* : to measure how a human would compute the action of a policy given a state at each environment step, we measure policy step inference time in seconds.

2. *Policy Size* : to measure how easily a human can read the entire policy, we measure its size in bytes. While this correlates with inference time for MLPs and linear models, tree-based policies may have large sizes but quick inference because they do not traverse all decision paths at each step.

As these measurements depend on many technical details (programming language, the compiler if any, the operating system, versions of libraries, the hardware it is executed on, etc), to ensure fair comparisons, we translate all student policies into a simple representation that mimics how a human being "reads" a policy. We call this process of standardizing policies language "unfolding". In Figure 9.1, 9.2, and 9.3, we present some unfolded policy programs. Other works have distilled neural networks into programs [103] or even directly learn programmatic policies [81] from scratch. However, those works directly consider programs as a policy class and could compare a generic program (not unfolded, with, e.g., while loops or array operations) to, e.g, a decision tree [99]. We will discuss later on the limitations of unfolding policies in the overall methodology.

## 9.2 Computing Baseline Policies

### 9.2.1 Setup

All the experiments presented next run on a dedicated cluster of Intel Xeon Gold 6130 (Skylake-SP), 2.10GHz, 2 CPUs/node, 16 cores/CPU with a timeout of 4 hours per experiment. Codes to reproduce our results are given in the supplementary material. In the future, we will open source a python library with all the tools of our methodology. Using Algorithm 10, we distill deep neural network expert policies into less complex policy classes.

**Policy classes** We consider four policy classes for our baselines. We choose those policy classes because there exist efficient algorithms to fit them with supervised data which is

```

1 import gymnasium as gym
2
3 env = gym.make("MountainCar")
4 s, _ = env.reset()
5 done = False
6 while not done:
7     y0 = 0.969*s[0]-30.830*s[1]
8     y1 = -0.205*s[0]+22.592*s[1]
9     y2 = -0.763*s[0]+8.237*s[1]
10    max_val = y0
11    action = 0
12    if y1 > max_val:
13        max_val = y1
14        action = 1
15    if y2 > max_val:
16        action = 2
17    s, r, terminated, truncated = env.step(action)
18    infos = env.step(action)
19    done = terminated or truncated

```

FIGURE 9.1 – Unfolded linear policy interacting with an environment.

```

1 def play(x):
2     h_layer_0_0 = 1.238*x[0]+0.971*x[1]
3     h_layer_0_0 = max(0, h_layer_0_0)
4     h_layer_0_1 = -1.221*x[0]+1.001
5     h_layer_0_1 = max(0, h_layer_0_1)
6     h_layer_1_0 = -0.109*h_layer_0_0
7     h_layer_1_0 = max(0, h_layer_1_0)
8     h_layer_1_1 = -3.024*h_layer_0_0
9     h_layer_1_1 = max(0, h_layer_1_1)
10    h_layer_2_0 = -1.790*h_layer_1_0
11    h_layer_2_0 = max(0, h_layer_2_0)
12    y_0 = h_layer_2_0
13    return [y_0]

```

FIGURE 9.2 – Tiny ReLU MLP for Pendulum.

```

1  def play(x):
2      if (x[3] - x[4]) <= -0.152:
3          if (x[2] - x[3]) <= 0.049:
4              return 3
5          else:
6              if x[5] <= 0.030:
7                  return 2
8              else:
9                  if (x[2] - x[4]) <= 0.049:
10                     return 3
11                 else:
12                     # Fire main engine
13                     (*@\colorbox{cyan}{return 1}@*)
14             else:
15                 if (x[0] - x[6]) <= -0.840:
16                     return 0
17                 else:
18                     if (x[0] - x[4]) <= -0.840:
19                         # Fire right engine
20                         (*@\colorbox{yellow}{return 1}@*)
21                     else:
22                         # Fire left engine
23                         (*@\colorbox{pink}{return 1}@*)

```

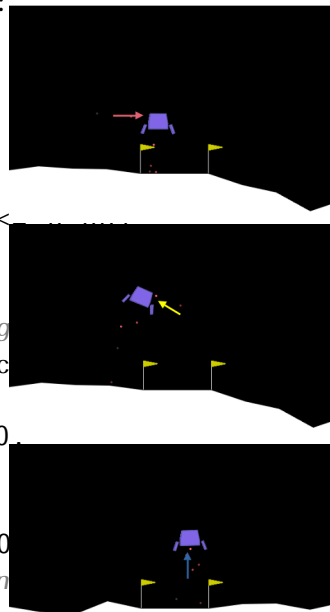


FIGURE 9.3 – An unfolded oblique tree policy’s actions obtaining 250 rewards on Lunar Lander.

Policy Class	Parameters	Training Algorithm
Linear Policies	Determined by state-action dimensions	Linear/Logistic Regression
Decision Trees	[4, 8, 16, 64, 128] nodes	CART (2× nodes maximum leaves)
Oblique Decision Trees	[4, 8, 16, 64, 128] nodes	CART (2× nodes maximum leaves)
ReLU MLPs	[2×2, 4×4, 8×8, 16×16] weigths	Adam optimization (500 iterations)

TABLEAU 9.1 – Summary of baseline policy classes parameters and fitting algorithms (used in Line 10).

a required step of imitation learning in Line 10. We consider linear policies that have been shown to be able to solve Mujoco tasks [64]. We fit linear policies to expert policies using simple linear (logistic) regressions with scikit-learn [77] default implementation. We also consider decision trees [**cart**] and oblique decision trees [**oblique**]. (Oblique) Decision trees are often considered the most interpretable model class in machine learning [**mythos**] and reinforcement learning [**IBMDP**, 6, 42, 70]. We train trees using the default CART [**cart**] implementation of scikit-learn with varying numbers of parameters (number of nodes in the tree). We also consider MLPs with ReLU activations [48] with varying number of parameters (total number of weights). This class of policy is often considered the least interpretable and is often used in deep reinforcement learning [47, 21, 49]. We train ReLU MLPs using the default scikit-learn implementation of Adam optimization [53] with 500 iterations. The 15 baseline policy classes that we consider are summarized in Appendix 9.1.

**Neural network experts** We do not train new deep reinforcement learning agents [71, 90, 47] but rather re-use ones available at the stables-baselines3 zoo [84]. Depending on the environments described next, we choose neural network policies from different deep reinforcement learning agents. Some may argue that during the imitation learning, ReLU MLPs baselines may obtain better performance because they are often from the same class as the expert they imitate unlike trees. But this is not of our concern as we do not benchmark the imitation learning algorithms. Furthermore, it is important to note that not all experts are compatible with all the variants of imitation learning Algorithm 10. Indeed, SAC experts [47] are not compatible with Q-Dagger [6] because it only works for continuous actions; and PPO experts, despite working with discrete actions do not compute a Q-function necessary for the re-weighting in Q-Dagger.

**Environments** We consider common environments in reinforcement learning research. We consider the classic control tasks from gymnasium [98], MuJoCo robots from [96], and Atari games from [7]. For Atari games, since the state space is frame pixels that can't

be interpreted, we use the object-centric version of the games from [27]. In Appendix 9.3 we give the list of environments we consider in our experiments with their state-action spaces as well as a cumulative reward threshold past which an environment is considered “solved”.

### 9.2.2 Ablation study of imitation learning

In this section, we present the results of the expert distillation into smaller policies. For each environment, we fit all the policy classes. To do so, we run different instantiations of Algorithm 10 multiple times with different total sample sizes. For each environment and each imitation learning variant, we summarize the number of times we fit all the baselines to an expert and which expert we use. The number of runs and imitation algorithm variants of Algorithm 10 are summarized in Appendix 9.4. After running the imitation learnings, we obtain roughly 40000 baseline policies (35000 for classic control, 5000 thousands for MuJoCo and 400 for OCArari). A dataset with all the baselines measurements is given in the supplementary material.

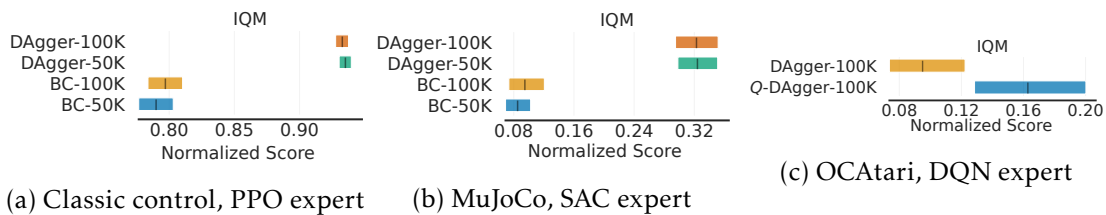


FIGURE 9.4 – Performance of imitation learning variants of Algorithm 10 on different environments. We plot the 95% stratified bootstrapped confidence intervals around the IQMs.

**What is the best imitation algorithm?** Even though the focus of our work is to evaluate trained policies, we still provide some insights on the best way to obtain interpretable policies from experts. Using the reinforcement learning evaluation library rliable [2], we plot on Figure 9.4 the interquartile means (IQM, an estimator of the mean robust to outliers) of the baseline policies cumulative rewards averaged over 100 episodes. For each imitation algorithm variant, we aggregate cumulative rewards over environments and policy classes. We normalize the baselines cumulative rewards between expert and random agent cumulative rewards.

The key observation is that for tested environments (Figures 9.4a, 9.4b), Behavior Cloning is not an efficient way to train baseline policies compared to DAGger. This is

probably because Behavior Cloning trains a student policy to match the expert’s actions on states visited by the expert while DAgger trains a student to take the expert’s actions on the states visited by the student [88]. An other observation is that the best performing imitation algorithms for MuJoCo (DAgger, Figure 9.4b) and OCArari (Q-Dagger, Figure 9.4c) obtain baselines that in average cannot match well the performances of the experts. However baseline policies almost always match the expert on simple tasks like classic control (Figure 9.4a).

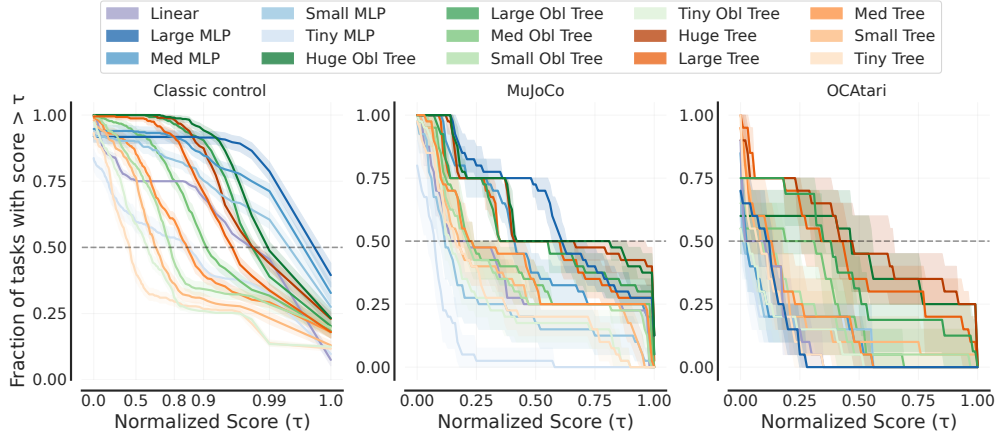


FIGURE 9.5 – Performance profiles of different policy classes on different environments.

**What is the best policy class in terms of reward?** We also wonder if there is a policy class that matches expert performances more often than others across environments. For that we plot performance profiles of the different policy classes obtained with a fixed expert and fixed imitation learning algorithm. In particular, for each environments group we use the baseline policies obtained from the best performing imitation learning algorithm from Figure 9.4. From Figure 9.5 we see that on classic control environments, MLPs tend to perform better than other classes while on OCArari games, trees tend to perform better than other classes. Now we move on to interpretability evaluation of our programmatic policies.

## 9.3 Measuring Policy Interpretability

### 9.3.1 From Policy to Program

In this section, we compute the step inference times, as well as the policy size for both the folded and unfolded variant of each policy obtained for classic control



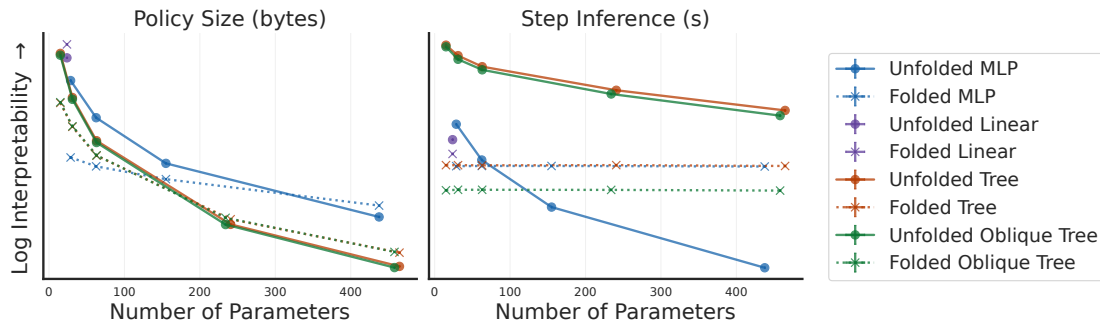


FIGURE 9.6 – Policies interpretability on classic control environments. We plot 95% stratified bootstrapped confidence intervals around means in both axes. In each sub-plot, interpretability is measured with either bytes or inference speed.

environments with DAGger-100K. To unfold policies, we convert them into Python programs formatted with PEP 8 (comparing other unfolding formats such as ONNX <https://github.com/onnx/onnx> is left to future work). We ensure that all policies operations are performed sequentially and compute the metrics for each policy on 100 episodes using the same CPUs.

**Is it necessary to unfold policies to compute interpretability metrics?** We see on Figure 9.6 that folded policies of the same class almost always give similar interpretability values (dotted lines) despite having very different number of parameters. Hence, measuring folded policies interpretability would contradict established results from user studies such as, e.g., trees of different sizes have different levels of interpretability [56].

**Is there a best policy class in terms of interpretability?** User studies from [36, 66, 102] show that decision trees are easier to understand than models involving mathematical equations like oblique trees, linear maps, and MLPs. However, [mythos] states that for a human wanting to have a global idea of the inference of a policy, a compact MLP can be more interpretable than a very deep decision tree. In Figure 9.6, we show that inference speed and memory size of programs help us capture those nuances : policy interpretability does not only depend on the policy class but also on the metric choice. Indeed, when we measure interpretability with inference times, we do observe that trees are more interpretable than MLPs. However, when measuring interpretability with policy size, we observe that MLPs can be more interpretable than trees for similar number of parameters. Because there seem to not be a more interpretable policy class across proxy metrics, we will keep studying both metrics at the same time.

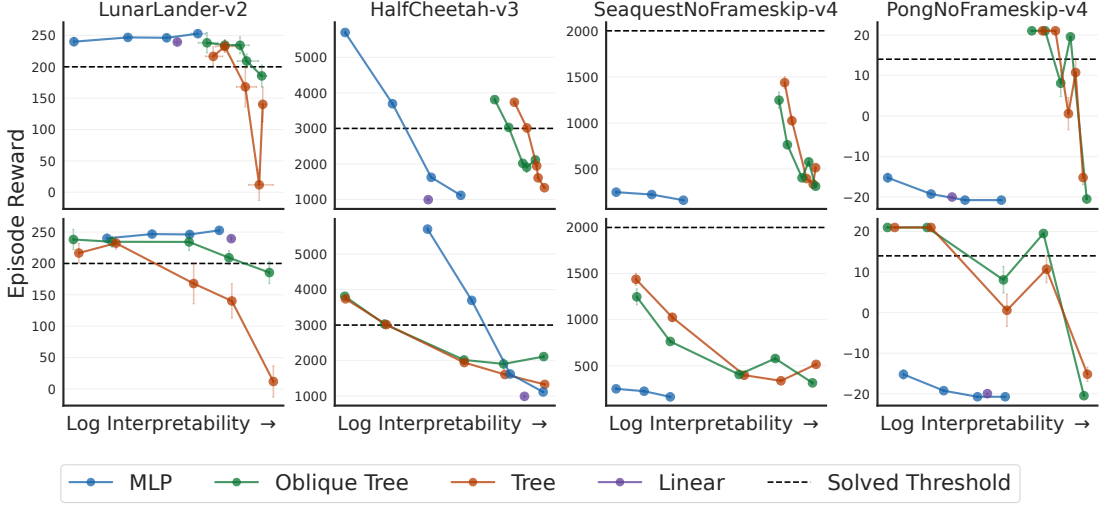


FIGURE 9.7 – Interpretability-Performance trade-offs. Top row, interpretability is measured with step inference times. Bottom row, the interpretability is measured with policy size. We plot 95% bootstrapped confidence intervals around means on both axes.

### 9.3.2 Interpretability-performance trade-offs

Now that we trained baseline policies and validated the proposed methodology, we use the latter to tackle open problems in interpretable reinforcement learning. For each environment, we fix the imitation learning algorithm and save the best baseline policy of each class in terms of episodic rewards after unfolding them. Each single Python policy is then **run again on the same dedicated CPU** for 100 new environment episodes (similarly to choosing a classifier with validation score and reporting the test score in the context of supervised learning).

**Is it possible to compute interpretable policies for high-dimensional environments?** [42] claim that computing an interpretable policy for high dimensional MDPs is difficult since it is similar to program synthesis which is known to be NP-hard [46]. Using our measures of interpretability, we can corroborate this claim. On Figure 9.7, we can indeed observe that some relatively interpretable policies can solve Pong (20 state dimensions) or HalfCheetah (17 state dimensions) while for very high-dimensional environments like Seaquest (180 state dimensions), no baseline can solve the game.

**For what environment are there good interpretable policies?** We fitted a random forest regressor [14] to predict the interpretability values of our baseline policies using

Environment Attributes	Importance for Step inference	Importance for Policy size
States dimension	<b>80.87</b>	<b>35.52</b>
Expert episodes lengths	11.39	9.28
Episode reward of random	2.26	4.75
Expert episode reward	1.51	16.80
Episode reward to solve	1.41	14.26
Actions dimension	1.41	2.02
Expert reward - Solve reward	1.15	17.37

TABLEAU 9.2 – Environment attributes importance to predict interpretability using either of our metrics.

environment attributes. In Table 9.2 we report the importance of each environment attribute when it comes to accurately predicting interpretability scores. We show that as hinted previously, the states dimensions of the environment is determining to predict the interpretability of good policies. Unsurprisingly, expert attributes also influence interpretability : for the environments where there is a positive large gap between expert and threshold rewards, the task could be considered easy and vice-versa.

**How does interpretability influence performance?** [64, 65] show the existence of linear and tree policies respectively that solve MuJoCo and continuous maze environments respectively; essentially showing that there exist environments for which policies more interpretable than deep neural networks can still compete performance-wise. Our evaluation indeed shows the existence of such environments. On Figure 9.7 we observe that on, e.g., LunarLander, increasing policy interpretability up to a certain point does not decrease reward. Actually, we can observe that for Pong a minimum level of interpretability is required to solve the game. Indeed, as stated in [35], optimizing interpretability can also be seen as regularizing the policy which can increase generalization capabilities. The key observation is that the policy class achieving the best interpretability-performance trade-off depends on the problem. Indeed, independent of the interpretability proxy metric, we see on Figure 9.7 that for LunarLander it is an MLP that achieves the best trade-off while for Pong it is a tree. Next, we compare our proxies for interpretability with another one; the verification time of policies used in [6, 5].

### 9.3.3 Verifying interpretable policies

[5] states that the cost of formally verifying properties of MLPs scales exponentially with the number of the parameters. Hence, they propose to measure interpretability of a

policy as the computations required to verify properties of actions given state subspaces, what they call local explainability queries [45]. Before [5], [6] also compared the time to formally verified properties of trees to the time to verify properties of MLPs to evaluate interpretability. In practice, this amounts to passing states and actions bounds and solving the SAT problem of finding a state in the state bounds for which the policy outputs an action in the action bounds. For example, for the LunarLander problem, a query could be to verify if when the y-position of the lander is below some threshold value, i.e, when the lander is close to the ground, there exists a state such that the tested policy would output the action of pushing towards the ground : if the solver outputs “SAT”, then there is a risk that the lander crashes.

Designing interesting queries covering all risks is an open problem, hence to evaluate the verification times of our baseline policies, we generate 500 random queries per environment by sampling state and action subspaces uniformly. Out of those queries we only report the verification times of “UNSAT” queries since to verify that, e.g., the lander does not crash we want the queries mentioned above to be “UNSAT”. We also only verify instances of ReLU MLPs using [110] for this experiment as verifying decision trees requires a different software [23] for which verification times would not be comparable.

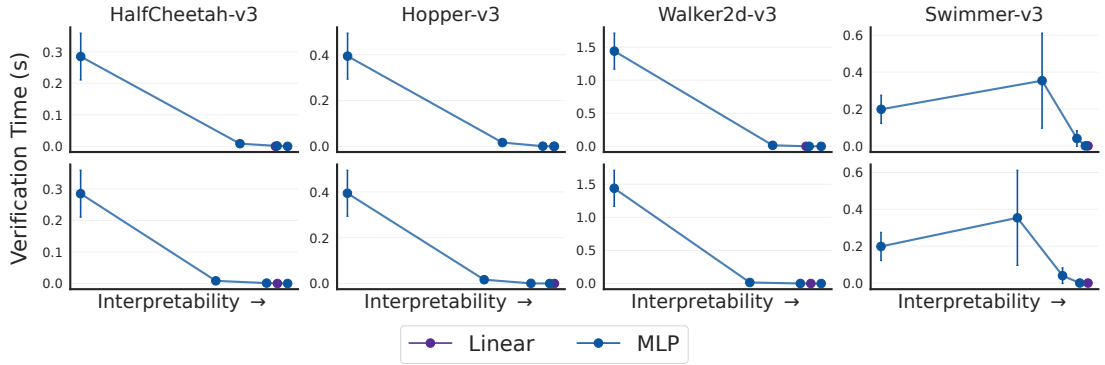


FIGURE 9.8 – Verification time as a function of policy interpretability. Top row, interpretability is measured with step inference times. Bottom row, the interpretability is measured with policy size. We plot 95% confidence intervals around means on both axes.

On Figure 9.8, we can observe that verification time decreases exponentially with MLP interpretability, both memory and inference speed, as shown in [5]. This is another good validation of our proposed methodology as well as a motivation to learn interpretable policies.

## 9.4 Experimental details

In this section we give all the experimental details necessary to reproduce our results.

Classic	MuJoCo	OCAtari
CartPole (4, 2, <b>490</b> )	Swimmer (8, 2, <b>300</b> )	Breakout (452, 4, <b>30</b> )
LunarLander (8, 4, <b>200</b> )	Walker2d (17, 6, <b>2000</b> )	Pong (20, 6, <b>14</b> )
LunarLanderContinuous (8, 2, <b>200</b> )	HalfCheetah (17, 6, <b>3000</b> )	SpaceInvaders (188, 6, <b>680</b> )
BipedalWalker (24, 4, <b>250</b> )	Hopper (11, 3, <b>2000</b> )	Seaquest (180, 18, <b>2000</b> )
MountainCar (2, 3, <b>90</b> )		
MountainCarContinuous (2, 1, <b>-110</b> )		
Acrobot (6, 3, <b>-100</b> )		
Pendulum (3, 1, <b>-400</b> )		

TABLEAU 9.3 – Summary of considered environments (dimensions of states and number or dimensions of actions, **reward thresholds**). The rewards thresholds are obtained from gymnasium [98]. For OCAtari environments, we choose the thresholds as the minimum between the DQN expert from [84] and the human scores. We also adapt subjectively some thresholds that we find too restrictive especially for MuJoCo (for example, the PPO expert from [84] has 2200 reward on Hopper while the default threshold was 3800).

## 9.5 All interpretability-performance trade-offs

In this appendix we provide the interpretability-performance trade-offs of all the tested environments. All the measures come from the experiment from Section 9.3.2.

Envs	BC 50K	BC 100K	Dagger 50K	Dagger 100K	Q 50K	Q-Dagger 100K
Classic	50 (PPO, DQN)	50 (PPO, DQN)	50 (PPO, DQN)	50 (PPO, DQN)	50 (DQN)	50 (DQN)
OCAtari	0	0	0	5 (DQN)	0	5 (DQN)
Mujoco	10 (SAC)	10 (SAC)	10 (SAC)	10 (SAC)	0	0

TABLEAU 9.4 – Repetitions of each imitation learning algorithm on each environment. We specify which deep reinforcement learning agent from the zoo [84] uses as experts in parentheses.

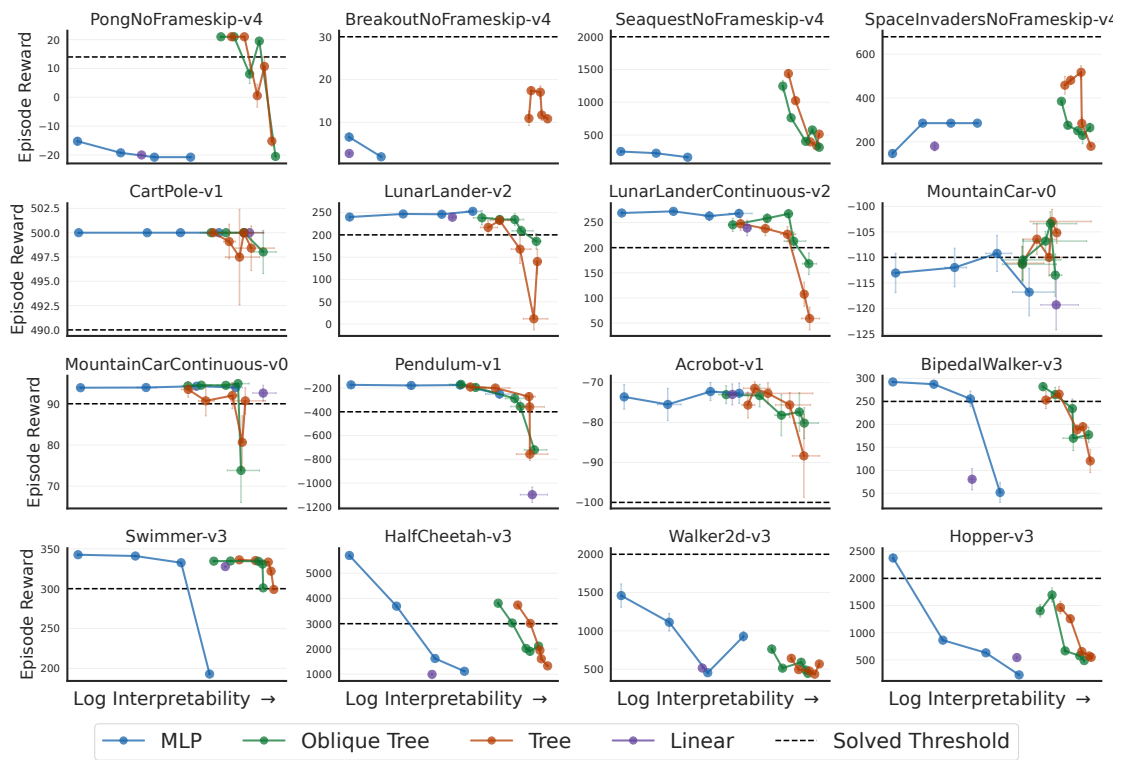


FIGURE 9.9 – Trade-off Cumulative Reward vs. Step Inference Time

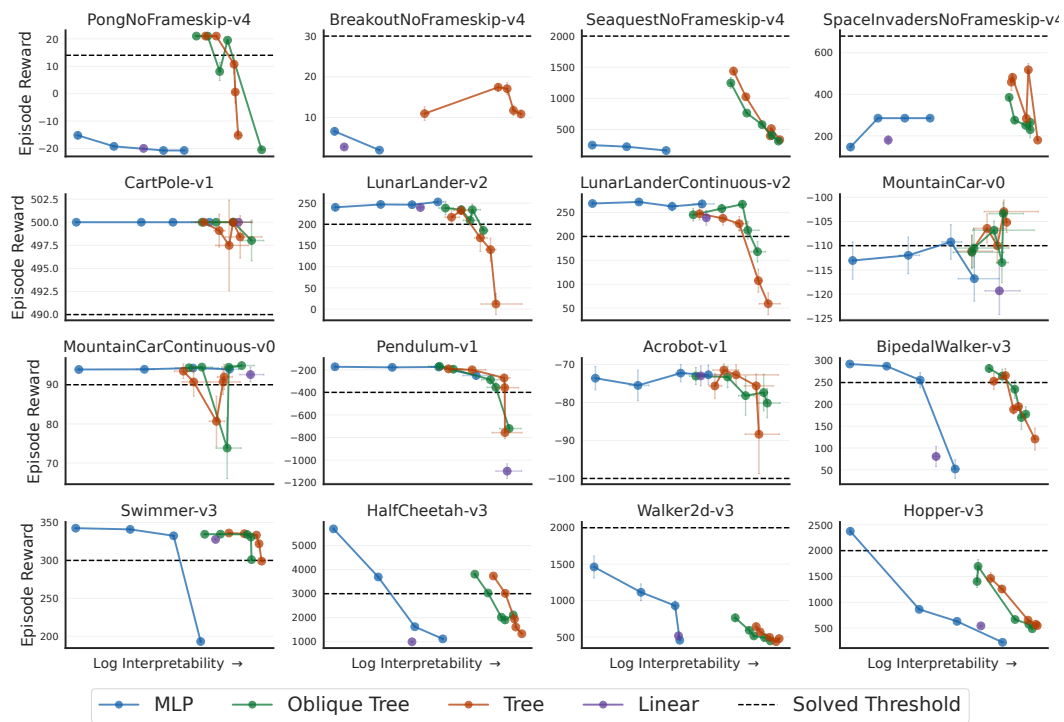


FIGURE 9.10 – Trade-off Cumulative Reward vs. Episode Inference Time

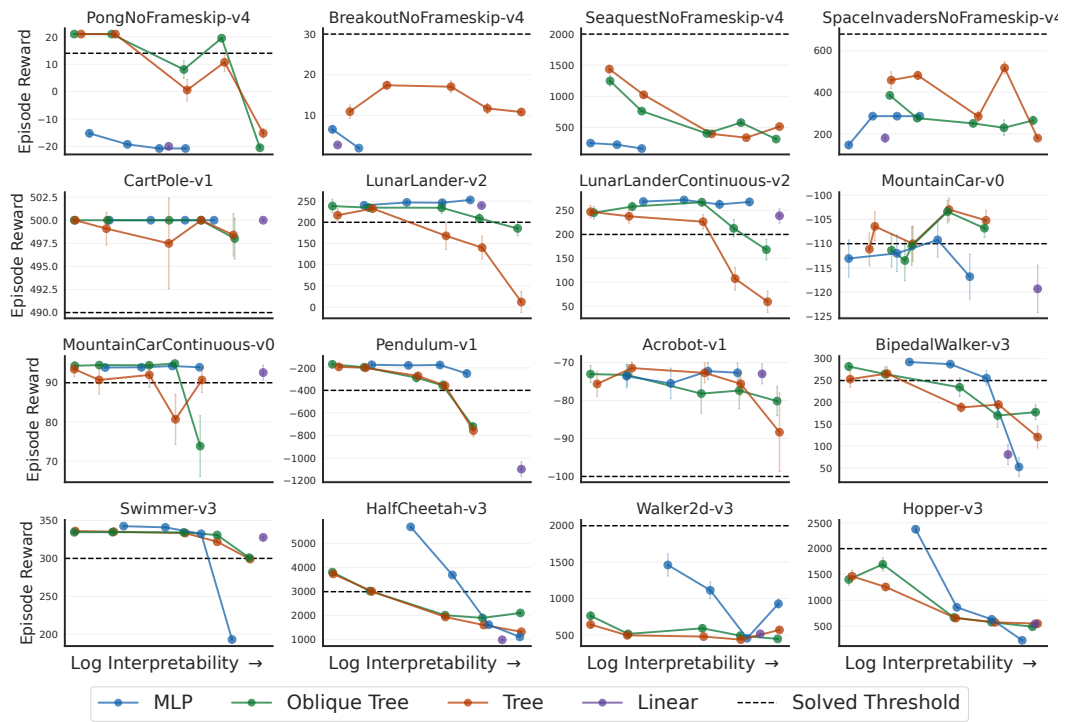


FIGURE 9.11 – Trade-off Cumulative Reward vs. Policy Size



# Conclusion Imitation

## 10.1 Limitations and conclusions

We have shown that our proposed methodology provides researchers with a sound way of evaluating policy interpretability. In particular, we have shown that unfolding policies in a common language such as Python is a key component of our methodology to ensure that interpretability depends on the policy complexity (c.f. Figure 9.6). Furthermore, we were able to show that the proxies we use for interpretability leads to similar conclusions from user studies of interpretability or from other empirical evaluations of interpretability (c.f. Figures 9.6, 9.7, and 9.8). Using the proposed methodology, we were able to illustrate the trade-offs between episodic reward and interpretability of policies from different classes (c.f. Figure 9.7) and showed the crucial need of our methodology as there is no better off policy class across tasks and metrics (c.f. Figures 9.5, 9.6, and 9.7).

A nice property of our methodology is that it is independent of the learning algorithm of the interpretable policy. We chose imitation learning but it could have been a random search in the policies parameters space [64]. Furthermore, there could be no limitation to use our methodology to evaluate the interpretability of arbitrary compositions of linear policies, trees and oblique trees, and MLPs, such as the hybrid policies from [92]. However, the unfolded version of policies with loops which lengths depend on the state would change between step, hence, the policy size metric value will change during episodes. This is not necessarily a strong limitation but would require more work on the unfolding procedures as well as on defining episodic metrics.

In the future, it would be interesting to compare episodic to averaged measures of

interpretability. Indeed, we additionally show in Appendix 9.10 the interpretability-performance trade-offs using the inference time summed over entire episodes as the measure of interpretability. Even though using episodic inference does not change the trade-offs compared to step inference time, it is important to discuss this nuance in future work since a key difference between supervised learning and reinforcement learning interpretability could be that human operators would read policies multiple times until the end of a decision process. Using episodic metrics for interpretability is not as straightforward as someone would think as for some MDPs, e.g. Acrobot, the episodes lengths depend on the policy. We also did not evaluate the role of sparsity in the interpretability of linear and MLP policies even though this could greatly influence the inference time. In the future it would be interesting to apply our methodologies to policies obtained with e.g. [93]. Moving away from evaluation, we also believe that our interpretable baselines can be used to train hierarchical agents [111] using our baselines as options. We hope that our methodology as well as the provided baselines will pave the way to a more rigorous science of interpretable reinforcement learning.

# Conclusion générale



# Bibliographie

- [1] Fernando ACERO et Zhibin LI. « Distilling Reinforcement Learning Policies for Interpretable Robot Locomotion : Gradient Boosting Machines and Symbolic Regression ». In : (2024). URL : <https://openreview.net/forum?id=fa3fjH3dEW>.
- [2] Rishabh AGARWAL et al. « Deep Reinforcement Learning at the Edge of the Statistical Precipice ». In : *Advances in Neural Information Processing Systems* (2021).
- [3] Sina AGHAEI, Andres GOMEZ et Phebe VAYANOS. « Learning Optimal Classification Trees : Strong Max-Flow Formulations ». In : (2020). arXiv : 2002.09142 [stat.ML].
- [4] Safa ALVER et Doina PRECUP. « An Attentive Approach for Building Partial Reasoning Agents from Pixels ». In : *Transactions on Machine Learning Research* (2024). ISSN : 2835-8856. URL : <https://openreview.net/forum?id=S3FUKFMRw8>.
- [5] Pablo BARCELÓ et al. « Model interpretability through the lens of computational complexity ». In : *Advances in neural information processing systems* (2020).
- [6] Osbert BASTANI, Yewen PU et Armando SOLAR-LEZAMA. « Verifiable Reinforcement Learning via Policy Extraction ». In : (2018).
- [7] Marc G. BELLEMARE et al. « The arcade learning environment : an evaluation platform for general agents ». In : *J. Artif. Int. Res.* 47.1 (mai 2013), p. 253-279. ISSN : 1076-9757.
- [8] James BERGSTRÄ, Daniel YAMINS et David COX. « Making a Science of Model Search : Hyperparameter Optimization in Hundreds of Dimensions for Vision Architectures ». In : *Proceedings of the 30th International Conference on Machine Learning*. Proceedings of Machine Learning Research 28.1 (17–19 Jun 2013). Sous la dir. de Sanjoy DASGUPTA et David McALLESTER, p. 115-123. URL : <https://proceedings.mlr.press/v28/bergstra13.html>.
- [9] Dimitris BERTSIMAS et Jack DUNN. « Optimal classification trees ». In : *Machine Learning* 106 (2017), p. 1039-1082.
- [10] Jock BLACKARD. « Covertypes ». In : (1998). DOI : <https://doi.org/10.24432/C50K5N>.
- [11] Guy BLANC et al. « Harnessing the power of choices in decision tree learning ». In : *Advances in Neural Information Processing Systems* 36 (2023), p. 80220-80232.

- [12] George BOOLE. *The Laws of Thought*. Walton, Maberly Macmillan et Co., 1854.
- [13] L BREIMAN et al. *Classification and Regression Trees*. Wadsworth, 1984.
- [14] Leo BREIMAN. « Random forests ». In : *Machine learning* 45 (2001), p. 5-32.
- [15] Marco BRESSAN et al. « A Theory of Interpretable Approximations ». In : *Proceedings of Thirty Seventh Conference on Learning Theory*. Proceedings of Machine Learning Research 247 (2024), p. 648-668.
- [16] Lars BUITINCK et al. « API design for machine learning software : experiences from the scikit-learn project ». In : *ECML PKDD Workshop : Languages for Data Mining and Machine Learning* (2013), p. 108-122.
- [17] Miguel A. CARREIRA-PERPINAN et Pooya TAVALLALI. « Alternating optimization of decision trees, with application to learning sparse oblique trees ». In : *Advances in Neural Information Processing Systems* 31 (2018). Sous la dir. de S. BENGIO et al. URL : [https://proceedings.neurips.cc/paper\\_files/paper/2018/file/185c29dc24325934ee377cfda20e414c-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2018/file/185c29dc24325934ee377cfda20e414c-Paper.pdf).
- [18] Miguel Á CARREIRA-PERPIÑÁN et Arman ZHARMAGAMBETOV. « Ensembles of Bagged TAO Trees Consistently Improve over Random Forests, AdaBoost and Gradient Boosting ». In : *Proceedings of the 2020 ACM-IMS on Foundations of Data Science Conference*. FODS '20 (2020), p. 35-46. DOI : 10.1145/3412815.3416882. URL : <https://doi.org/10.1145/3412815.3416882>.
- [19] Ayman CHAOUKI, Jesse READ et Albert BIFET. « Branches : A Fast Dynamic Programming and Branch & Bound Algorithm for Optimal Decision Trees ». In : (2024). arXiv : 2406.02175 [cs.LG]. URL : <https://arxiv.org/abs/2406.02175>.
- [20] Tianqi CHEN et Carlos GUESTRIN. « XGBoost : A Scalable Tree Boosting System ». In : *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* (2016), p. 785-794.
- [21] Xinyue CHEN et al. « Randomized Ensembled Double Q-Learning : Learning Fast Without a Model ». In : (2021). URL : <https://openreview.net/forum?id=AY8zfZm0tDd>.
- [22] Vinícius G COSTA et Carlos E PEDREIRA. « Recent advances in decision trees : An updated survey ». In : *Artificial Intelligence Review* 56 (2023), p. 4765-4800.
- [23] Leonardo DE MOURA et Nikolaj BJØRNER. « Z3 : an efficient SMT solver ». In : *TACAS'08/ETAPS'08* (2008), p. 337-340.
- [24] Jonas DEGRAVE et al. « Magnetic control of tokamak plasmas through deep reinforcement learning ». In : *Nature* 602.7897 (2022), p. 414-419.
- [25] Quentin DELFOSSE et al. « Interpretable and Explainable Logical Policies via Neurally Guided Symbolic Abstraction ». In : *Advances in Neural Information Processing (NeurIPS)* (2023).

- [26] Quentin DELFOSSE et al. « Interpretable Concept Bottlenecks to Align Reinforcement Learning Agents ». In : (2024). URL : <https://openreview.net/forum?id=ZC0PSk6Mc6>.
- [27] Quentin DELFOSSE et al. « OCArari : Object-Centric Atari 2600 Reinforcement Learning Environments ». In : *Reinforcement Learning Journal* 1 (2024), p. 400-449.
- [28] Emir DEMIROVIC et al. « MurTree : Optimal Decision Trees via Dynamic Programming and Search ». In : *Journal of Machine Learning Research* 23.26 (2022), p. 1-47. URL : <http://jmlr.org/papers/v23/20-520.html>.
- [29] Emir DEMIROVIĆ, Emmanuel HEBRARD et Louis JEAN. « Blossom : an Anytime Algorithm for Computing Optimal Decision Trees ». In : *Proceedings of the 40th International Conference on Machine Learning*. Proceedings of Machine Learning Research 202 (23-29 Jul 2023). Sous la dir. d'Andreas KRAUSE et al., p. 7533-7562. URL : <https://proceedings.mlr.press/v202/demirovic23a.html>.
- [30] Jacob DEVLIN et al. « Bert : Pre-training of deep bidirectional transformers for language understanding ». In : *Proceedings of the 2019 conference of the North American chapter of the association for computational linguistics : human language technologies, volume 1 (long and short papers)*. 2019, p. 4171-4186.
- [31] Finale DOSHI-VELEZ et Been KIM. « Towards A Rigorous Science of Interpretable Machine Learning ». In : (2017). arXiv : 1702.08608 [stat.ML]. URL : <https://arxiv.org/abs/1702.08608>.
- [32] Gabriel DULAC-ARNOLD et al. « Datum-Wise Classification : A Sequential Approach to Sparsity ». In : *Machine Learning and Knowledge Discovery in Databases* (2011), p. 375-390. ISSN : 1611-3349. DOI : 10.1007/978-3-642-23780-5\_34. URL : [http://dx.doi.org/10.1007/978-3-642-23780-5\\_34](http://dx.doi.org/10.1007/978-3-642-23780-5_34).
- [33] Jan-Niklas ECKARDT et al. « Reinforcement learning for precision oncology ». In : *Cancers* 13.18 (2021), p. 4624.
- [34] Floriana ESPOSITO et al. « A comparative analysis of methods for pruning decision trees ». In : *IEEE transactions on pattern analysis and machine intelligence* 19.5 (1997), p. 476-491.
- [35] Alex A. FREITAS. « Comprehensible classification models : a position paper ». In : *SIGKDD Explor. Newsl.* 15.1 (mars 2014), p. 1-10. ISSN : 1931-0145. DOI : 10.1145/2594473.2594475. URL : <https://doi.org/10.1145/2594473.2594475>.
- [36] Alex A. FREITAS, Daniela C. WIESER et Rolf APWEILER. « On the Importance of Comprehensible Classification Models for Protein Function Prediction ». In : *IEEE/ACM Trans. Comput. Biol. Bioinformatics* 7.1 (jan. 2010), p. 172-182. ISSN : 1545-5963. DOI : 10.1109/TCBB.2008.47. URL : <https://doi.org/10.1109/TCBB.2008.47>.

- [37] Yoav FREUND et Robert E SCHAPIRE. « A Decision-Theoretic Generalization of On-Line Learning and an Application to Boosting ». In : *Journal of Computer and System Sciences* 55.1 (1997), p. 119-139. ISSN : 0022-0000. DOI : <https://doi.org/10.1006/jcss.1997.1504>. URL : <https://www.sciencedirect.com/science/article/pii/S002200009791504X>.
- [38] Jerome H. FRIEDMAN. « Greedy Function Approximation : A Gradient Boosting Machine ». In : *The Annals of Statistics* 29.5 (2001), p. 1189-1232.
- [39] Jerome H. FRIEDMAN. « Stochastic gradient boosting ». In : *Comput. Stat. Data Anal.* 38.4 (2002), p. 367-378.
- [40] Abhinav GARLAPATI et al. « A Reinforcement Learning Approach to Online Learning of Decision Trees ». In : (2015). arXiv : 1507.06923 [cs.LG]. URL : <https://arxiv.org/abs/1507.06923>.
- [41] Romain GAUTRON. « FApprentissage par renforcement pour l'aide à la conduite des cultures des petits agriculteurs des pays du Sud : vers la maîtrise des risques. » Thèse de doct. Montpellier SupAgro, 2022.
- [42] Claire GLANOIS et al. « A survey on interpretable reinforcement learning ». In : *Machine Learning* (2024), p. 1-44.
- [43] Yury GORISHNIY et al. « Revisiting deep learning models for tabular data ». In : *Proceedings of the 35th International Conference on Neural Information Processing Systems* (2024).
- [44] LÉO GRINSZTAJN, Edouard OYALLON et Gaël VAROQUAUX. « Why do tree-based models still outperform deep learning on typical tabular data? » In : *Advances in neural information processing systems* 35 (2022), p. 507-520.
- [45] Riccardo GUIDOTTI et al. « A Survey of Methods for Explaining Black Box Models ». In : *ACM Comput. Surv.* 51.5 (août 2018). ISSN : 0360-0300. DOI : 10.1145/3236009. URL : <https://doi.org/10.1145/3236009>.
- [46] Sumit GULWANI, Oleksandr POLOZOV, Rishabh SINGH et al. « Program synthesis ». In : *Foundations and Trends® in Programming Languages* 4.1-2 (2017), p. 1-119.
- [47] Tuomas HAARNOJA et al. « Soft Actor-Critic : Off-Policy Maximum Entropy Deep Reinforcement Learning with a Stochastic Actor ». In : *Proceedings of the 35th International Conference on Machine Learning*. Sous la dir. de Jennifer DY et Andreas KRAUSE. T. 80. Proceedings of Machine Learning Research. PMLR, oct. 2018, p. 1861-1870. URL : <https://proceedings.mlr.press/v80/haarnoja18b.html>.
- [48] Kaiming HE et al. « Delving deep into rectifiers : Surpassing human-level performance on imagenet classification ». In : (2015), p. 1026-1034.
- [49] Takuya HIRAOKA et al. « Dropout Q-Functions for Doubly Efficient Reinforcement Learning ». In : (2022). URL : <https://openreview.net/forum?id=xCVJMsPv3RT>.



- [50] Laurent HYAFIL et Ronald L. RIVEST. « Constructing optimal binary decision trees is NP-complete ». In : *Information Processing Letters* 5.1 (1976), p. 15-17. ISSN : 0020-0190. DOI : [https://doi.org/10.1016/0020-0190\(76\)90095-8](https://doi.org/10.1016/0020-0190(76)90095-8). URL : <https://www.sciencedirect.com/science/article/pii/0020019076900958>.
- [51] Rasul KAIRGELDIN et Miguel Á. CARREIRA-PERPIÑÁN. « Bivariate Decision Trees : Smaller, Interpretable, More Accurate ». In : *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*. KDD '24 (2024), p. 1336-1347. DOI : 10.1145/3637528.3671903. URL : <https://doi.org/10.1145/3637528.3671903>.
- [52] Guolin KE et al. « Lightgbm : A highly efficient gradient boosting decision tree ». In : *Advances in neural information processing systems* 30 (2017), p. 3146-3154.
- [53] Diederik P. KINGMA et Jimmy BA. « Adam : A Method for Stochastic Optimization ». In : (2015).
- [54] Donald Ervin KNUTH. « Finite semifields and projective planes ». Thèse de doct. California Institute of Technology, 1963.
- [55] Brent KOMER, James BERGSTRA et Chris ELIASMITH. « Hyperopt-Sklearn : Automatic Hyperparameter Configuration for Scikit-Learn ». In : *Proceedings of the 13th Python in Science Conference* (2014). Sous la dir. de Stéfan van der WALT et James BERGSTRA, p. 32-37. DOI : 10.25080/Majora-14bd3278-006.
- [56] Nada LAVRAČ. « Selected techniques for data mining in medicine ». In : *Artificial Intelligence in Medicine* 16.1 (1999). Data Mining Techniques and Applications in Medicine, p. 3-23. ISSN : 0933-3657. DOI : [https://doi.org/10.1016/S0933-3657\(98\)00062-1](https://doi.org/10.1016/S0933-3657(98)00062-1). URL : <https://www.sciencedirect.com/science/article/pii/S0933365798000621>.
- [57] Yann LECUN et al. « Backpropagation applied to handwritten zip code recognition ». In : *Neural computation* 1.4 (1989), p. 541-551.
- [58] Edouard LEURENT. « Safe and Efficient Reinforcement Learning for Behavioural Planning in Autonomous Driving ». Thèse de doct. Université de Lille, 2020.
- [59] Jimmy LIN et al. « Generalized and scalable optimal sparse decision trees ». In : *International Conference on Machine Learning* (2020), p. 6150-6160.
- [60] Jacobus van der LINDEN, Mathijs de WEERDT et Emir DEMIROVIĆ. « Necessary and Sufficient Conditions for Optimal Decision Trees using Dynamic Programming ». In : *Advances in Neural Information Processing Systems* 36 (2023). Sous la dir. d'A. OH et al., p. 9173-9212.
- [61] Jacobus G. M. van der LINDEN et al. « Optimal or Greedy Decision Trees? Revisiting their Objectives, Tuning, and Performance ». In : (2024). arXiv : 2409.12788 [cs.LG]. URL : <https://arxiv.org/abs/2409.12788>.

- [62] Zachary C. LIPTON. « The Mythos of Model Interpretability : In machine learning, the concept of interpretability is both important and slippery. » In : *Queue* 16.3 (2018), p. 31-57.
- [63] Wei-Yin LOH. « Fifty years of classification and regression trees ». In : *International Statistical Review* 82.3 (2014), p. 329-348.
- [64] Horia MANIA, Aurelia GUY et Benjamin RECHT. « Simple random search of static linear policies is competitive for reinforcement learning ». In : 31 (2018). Sous la dir. de S. BENGIO et al. URL : [https://proceedings.neurips.cc/paper\\_files/paper/2018/file/7634ea65a4e6d9041cf7de18e334a-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2018/file/7634ea65a4e6d9041cf7de18e334a-Paper.pdf).
- [65] Yishay MANSOUR, Michal MOSHKOVITZ et Cynthia RUDIN. *There is no Accuracy-Interpretability Tradeoff in Reinforcement Learning for Mazes*. 2022. arXiv : 2206.04266 [cs.LG]. URL : <https://arxiv.org/abs/2206.04266>.
- [66] David MARTENS et al. « Performance of classification models from a user perspective ». In : *Decision Support Systems* 51.4 (2011). Recent Advances in Data, Text, and Media Mining & Information Issues in Supply Chain and in Service System Design, p. 782-793. ISSN : 0167-9236. DOI : <https://doi.org/10.1016/j.dss.2011.01.013>. URL : <https://www.sciencedirect.com/science/article/pii/S016792361100042X>.
- [67] Sascha MARTON et al. « Mitigating Information Loss in Tree-Based Reinforcement Learning via Direct Optimization ». In : (2025). URL : <https://openreview.net/forum?id=qpXctF2aLZ>.
- [68] Rahul MAZUMDER, Xiang MENG et Haoyue WANG. « Quant-BnB : A Scalable Branch-and-Bound Method for Optimal Decision Trees with Continuous Features ». In : *Proceedings of the 39th International Conference on Machine Learning*. Proceedings of Machine Learning Research 162 (17–23 Jul 2022). Sous la dir. de Kamalika CHAUDHURI et al., p. 15255-15277. URL : <https://proceedings.mlr.press/v162/mazumder22a.html>.
- [69] Ameet Talwalkar MEHRYAR MOHRI Afshin Rostamizadeh. *Foundations of Machine Learning*. MIT Press, 2012.
- [70] Stephanie MILANI et al. « Explainable Reinforcement Learning : A Survey and Comparative Review ». In : *ACM Comput. Surv.* 56.7 (avr. 2024). ISSN : 0360-0300. DOI : 10.1145/3616864. URL : <https://doi.org/10.1145/3616864>.
- [71] Volodymyr MNIH et al. « Human-level control through deep reinforcement learning ». In : *nature* 518.7540 (2015), p. 529-533.
- [72] W Nor Haizan W MOHAMED, Mohd Najib Mohd SALLEH et Abdul Halim OMAR. « A comparative study of reduced error pruning method in decision tree algorithms ». In : *2012 IEEE International conference on control system, computing and engineering* (2012), p. 392-397.

- [73] Sreerama MURTHY et Steven SALZBERG. « Decision tree induction : how effective is the greedy heuristic? » In : *Proceedings of the First International Conference on Knowledge Discovery and Data Mining* (1995), p. 222-227.
- [74] Sreerama MURTHY et Steven SALZBERG. « Lookahead and Pathology in Decision Tree Induction ». In : *Proceedings of the 14th International Joint Conference on Artificial Intelligence - Volume 2. IJCAI'95* (1995), p. 1025-1031.
- [75] Sreerama K MURTHY, Simon KASIF et Steven SALZBERG. « A system for induction of oblique decision trees ». In : *Journal of artificial intelligence research* 2 (1994), p. 1-32.
- [76] Mohammad NOROUZI et al. « Efficient Non-greedy Optimization of Decision Trees ». In : *Advances in Neural Information Processing Systems* 28 (2015). Sous la dir. de C. CORTES et al. URL : [https://proceedings.neurips.cc/paper\\_files/paper/2015/file/1579779b98ce9edb98dd85606f2c119d-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2015/file/1579779b98ce9edb98dd85606f2c119d-Paper.pdf).
- [77] F. PEDREGOSA et al. « Scikit-learn : Machine Learning in Python ». In : *Journal of Machine Learning Research* 12 (2011), p. 2825-2830.
- [78] Dean A POMERLEAU. « Alvin : An autonomous land vehicle in a neural network ». In : *Advances in neural information processing systems* 1 (1988).
- [79] Liudmila PROKHORENKOVA et al. « CatBoost : unbiased boosting with categorical features ». In : *Proceedings of the 32nd International Conference on Neural Information Processing Systems. NIPS'18* (2018), p. 6639-6649.
- [80] Martin L. PUTERMAN. *Markov Decision Processes : Discrete Stochastic Dynamic Programming*. John Wiley & Sons, 1994.
- [81] Wenjie QIU et He ZHU. « Programmatic Reinforcement Learning without Oracles ». In : (2022). URL : <https://openreview.net/forum?id=6Tk2noBdvxt>.
- [82] J ROSS QUINLAN. « C4. 5 : Programs for machine learning ». In : *Morgan Kaufmann google schola* 2 (1993), p. 203-228.
- [83] J. R. QUINLAN. « Induction of Decision Trees ». In : *Mach. Learn.* 1.1 (1986), p. 81-106.
- [84] Antonin RAFFIN. *RL Baselines3 Zoo*. GitHub, 2020.
- [85] Antonin RAFFIN et al. « Stable-Baselines3 : Reliable Reinforcement Learning Implementations ». In : *Journal of Machine Learning Research* 22.268 (2021), p. 1-8.
- [86] Marco Tulio RIBEIRO, Sameer SINGH et Carlos GUESTRIN. « "Why Should I Trust You?" : Explaining the Predictions of Any Classifier ». In : *KDD '16* (2016), p. 1135-1144. DOI : 10.1145/2939672.2939778. URL : <https://doi.org/10.1145/2939672.2939778>.
- [87] Frank ROSENBLATT. « The perceptron : a probabilistic model for information storage and organization in the brain. » In : *Psychological review* 65.6 (1958), p. 386.

- [88] Stéphane Ross, Geoffrey J. GORDON et J. Andrew BAGNELL. « A Reduction of Imitation Learning and Structured Prediction to No-Regret Online Learning ». In : (2010).
- [89] Patrick SAUX et al. « Development and validation of an interpretable machine learning-based calculator for predicting 5-year weight trajectories after bariatric surgery : a multinational retrospective cohort SOPHIA study ». In : *The Lancet Digital Health* (août 2023). DOI : 10.1016/S2589-7500(23)00135-8. URL : <https://hal.science/hal-04192198>.
- [90] John SCHULMAN et al. « Proximal policy optimization algorithms ». In : *arXiv preprint arXiv:1707.06347* (2017).
- [91] Yijun SHAO et al. « Shedding Light on the Black Box : Explaining Deep Neural Network Prediction of Clinical Outcomes ». In : *J. Med. Syst.* 45.1 (jan. 2021). ISSN : 0148-5598. DOI : 10.1007/s10916-020-01701-8. URL : <https://doi.org/10.1007/s10916-020-01701-8>.
- [92] Hikaru SHINDO et al. « BlendRL : A Framework for Merging Symbolic and Neural Policy Learning ». In : *arXiv* (2025).
- [93] Anna SOLIGO, Pietro FERRARO et David BOYLE. *Induced Modularity and Community Detection for Functionally Interpretable Reinforcement Learning*. 2025. arXiv : 2501.17077 [cs.LG]. URL : <https://arxiv.org/abs/2501.17077>.
- [94] Gowthami SOMEPALLI et al. « SAINT : Improved Neural Networks for Tabular Data via Row Attention and Contrastive Pre-Training ». In : (2021). arXiv : 2106.01342 [cs.LG]. URL : <https://arxiv.org/abs/2106.01342>.
- [95] Richard S. SUTTON et Andrew G. BARTO. *Reinforcement Learning : An Introduction*. Cambridge, MA : The MIT Press, 1998.
- [96] Emanuel TODOROV, Tom EREZ et Yuval TASSA. « MuJoCo : A physics engine for model-based control. » In : (2012), p. 5026-5033.
- [97] Nicholay TOPIN et al. « Iterative bounding mdps : Learning interpretable policies via non-interpretable methods ». In : *Proceedings of the AAAI Conference on Artificial Intelligence* 35 (2021), p. 9923-9931.
- [98] Mark TOWERS et al. « Gymnasium : A Standard Interface for Reinforcement Learning Environments ». In : *arXiv preprint arXiv:2407.17032* (2024).
- [99] Dweep TRIVEDI et al. « Learning to Synthesize Programs as Interpretable and Generalizable Policies ». In : (2021). Sous la dir. d'A. BEYGELZIMER et al. URL : <https://openreview.net/forum?id=wP9twkexC3V>.
- [100] Alan TURING. « Computing Machinery and Intelligence ». In : *Mind* (1950).
- [101] Joaquin VANSCHOREN et al. « OpenML : networked science in machine learning ». In : *SIGKDD Explor. Newsl.* 15.2 (juin 2014), p. 49-60. ISSN : 1931-0145. DOI : 10.1145/2641190.2641198. URL : <https://doi.org/10.1145/2641190.2641198>.

- [102] Wouter VERBEKE et al. « Building comprehensible customer churn prediction models with advanced rule induction techniques ». In : *Expert Systems with Applications* 38.3 (2011), p. 2354-2364. ISSN : 0957-4174. DOI : <https://doi.org/10.1016/j.eswa.2010.08.023>. URL : <https://www.sciencedirect.com/science/article/pii/S0957417410008067>.
- [103] Abhinav VERMA et al. « Programmatically interpretable reinforcement learning ». In : (2018), p. 5045-5054.
- [104] Sicco VERWER et Yingqian ZHANG. « Learning decision trees with flexible constraints and objectives using integer optimization ». In : *Integration of AI and OR Techniques in Constraint Programming : 14th International Conference, CPAIOR 2017, Padua, Italy, June 5-8, 2017, Proceedings* 14 (2017), p. 94-103.
- [105] Sicco VERWER et Yingqian ZHANG. « Learning optimal classification trees using a binary linear program formulation ». In : *Proceedings of the AAAI conference on artificial intelligence* 33 (2019), p. 1625-1632.
- [106] Daniël Vos et Sicco VERWER. « Optimal decision tree policies for Markov decision processes ». In : *Proceedings of the Thirty-Second International Joint Conference on Artificial Intelligence. IJCAI '23* (2023). DOI : 10.24963/ijcai.2023/606. URL : <https://doi.org/10.24963/ijcai.2023/606>.
- [107] Daniël Vos et Sicco VERWER. « Optimizing Interpretable Decision Tree Policies for Reinforcement Learning ». In : (2024). arXiv : 2408.11632 [cs.LG]. URL : <https://arxiv.org/abs/2408.11632>.
- [108] Maxime WABARTHA et Joelle PINEAU. « Piecewise Linear Parametrization of Policies : Towards Interpretable Deep Reinforcement Learning ». In : (2024). URL : <https://openreview.net/forum?id=h0MVq57Ce0>.
- [109] Daniel WHITESON. « HIGGS ». In : (2014). DOI : <https://doi.org/10.24432/C5V312>.
- [110] Haoze WU et al. *Marabou 2.0 : A Versatile Formal Analyzer of Neural Networks*. 2024. arXiv : 2401.14461 [cs.AI]. URL : <https://arxiv.org/abs/2401.14461>.
- [111] Jesse ZHANG, Haonan YU et Wei XU. « Hierarchical Reinforcement Learning by Discovering Intrinsic Options ». In : (2021). URL : <https://openreview.net/forum?id=r-gPPHEjpmw>.
- [112] Arman ZHARMAGAMBETOV, Magzhan GABIDOLLA et Miguel Ê. CARREIRA-PERPIÑÁN. « Improved Boosted Regression Forests Through Non-Greedy Tree Optimization ». In : *2021 International Joint Conference on Neural Networks (IJCNN)* (2021), p. 1-8. DOI : 10.1109/IJCNN52387.2021.9534446.
- [113] Arman ZHARMAGAMBETOV et al. « Non-Greedy Algorithms for Decision Tree Optimization : An Experimental Comparison ». In : *2021 International Joint Conference on Neural Networks (IJCNN)* (2021), p. 1-8. DOI : 10.1109/IJCNN52387.2021.9533597.



## Programmes informatiques

Les listings suivants sont au cœur de notre travail.

Listing A.1 – Il est l’heure

```

1  #include <stdio.h>
2  int heures, minutes, secondes;
3
4  /*****
5   *
6   *          print_heure
7   *
8   *   But:
9   *   Imprime l'heure *****/
10 /* *****/
11 /*...Interface: *****/
12 /*...Utilise les variables globales *****/
13 /*...heures, minutes, secondes *****/
14 /* *****/
15 /*****
16
17 void print_heure(void)
18 {
19     printf("Il est %d heure", heures);
20     if (heures > 1) printf("s");
21     printf(" %d minute", minutes);
22     if (minutes > 1) printf("s");
23     printf(" %d seconde", secondes);
24     if (secondes > 1) printf("s");
25     printf("\n");
26 }
```

## Listing A.2 – Factorielle

```
1 | int factorielle(int n)
2 | {
3 |     if (n > 2) return n * factorielle(n - 1);
4 |     return n;
5 | }
```



# Appendix I

## B.1 Tree value computations

**Depth-0 decision tree :** has only one leaf node that takes a single base action indefinitely. For this type of tree the best reward achievable is to take actions that maximize the probability of reaching the objective  $\rightarrow$  or  $\downarrow$ . In that case the objective value of such tree is : In the goal state  $G = (1, 0)$ , the value of the depth-0 tree  $\mathcal{T}_0$  is :

$$\begin{aligned} V_G^{\mathcal{T}_0} &= 1 + \gamma + \gamma^2 + \dots \\ &= \sum_{t=0}^{\infty} \gamma^t \\ &= \frac{1}{1 - \gamma} \end{aligned}$$

In the state  $(0, 0)$  when the policy repeats going right respectively in the state  $(0, 1)$  when the policy repeats going down, the value is :

$$\begin{aligned} V_{S_0}^{\mathcal{T}_0} &= 0 + \gamma V_g^{\mathcal{T}_0} \\ &= \gamma V_G^{\mathcal{T}_0} \end{aligned}$$

In the other states the policy never gets positive rewards ;  $V_{S_1}^{T_0} = V_{S_2}^{T_0} = 0$ . Hence :

$$\begin{aligned}
 J(T_0) &= \frac{1}{4}V_G^{T_0} + \frac{1}{4}V_{S_0}^{T_0} + \frac{1}{4}V_{S_1}^{T_0} + \frac{1}{4}V_{S_2}^{T_0} \\
 &= \frac{1}{4}V_G^{T_0} + \frac{1}{4}\gamma V_G^{T_0} + 0 + 0 \\
 &= \frac{1}{4}\frac{1}{1-\gamma} + \frac{1}{4}\gamma\frac{1}{1-\gamma} \\
 &= \frac{1+\gamma}{4(1-\gamma)}
 \end{aligned}$$

**Unbalanced depth-2 decision tree :** the unbalanced depth-2 decision tree takes an information gathering action  $x \leq 0.5$  then either takes the  $\downarrow$  action or takes a second information  $y \leq 0.5$  followed by  $\rightarrow$  or  $\downarrow$ . In states  $G$  and  $S_2$ , the value of the unbalanced tree is the same as for the depth-1 tree. In states  $S_0$  and  $S_1$ , the policy takes two information gathering actions before taking a base action and so on :

$$V_{S_0}^{T_u} = \zeta + \gamma\zeta + \gamma^2 0 + \gamma^3 V_G^{T_1}$$

$$\begin{aligned}
 V_{S_1}^{T_u} &= \zeta + \gamma\zeta + \gamma^2 0 + \gamma^3 V_{S_0}^{T_u} \\
 &= \zeta + \gamma\zeta + \gamma^2 0 + \gamma^3(\zeta + \gamma\zeta + \gamma^2 0 + \gamma^3 V_G^{T_1}) \\
 &= \zeta + \gamma\zeta + \gamma^3\zeta + \gamma^4\zeta + \gamma^6 V_G^{T_1}
 \end{aligned}$$

We get :

$$\begin{aligned}
 J(T_u) &= \frac{1}{4}V_G^{T_u} + \frac{1}{4}V_{S_0}^{T_u} + \frac{1}{4}V_{S_1}^{T_u} + \frac{1}{4}V_{S_2}^{T_u} \\
 &= \frac{1}{4}V_G^{T_1} + \frac{1}{4}(\zeta + \gamma\zeta + \gamma^3 V_G^{T_1}) + \frac{1}{4}(\zeta + \gamma\zeta + \gamma^3\zeta + \gamma^4\zeta + \gamma^6 V_G^{T_1}) + \frac{1}{4}V_{S_2}^{T_1} \\
 &= \frac{1}{4}\left(\frac{\zeta + \gamma}{1-\gamma^2}\right) + \frac{1}{4}\left(\frac{\gamma\zeta + \gamma^4 + \zeta - \gamma^2\zeta}{1-\gamma^2}\right) + \frac{1}{4}(\zeta + \gamma\zeta + \gamma^3\zeta + \gamma^4\zeta + \gamma^6 V_G^{T_1}) + \frac{1}{4}V_{S_2}^{T_1} \\
 &= \frac{1}{4}\left(\frac{\zeta + \gamma}{1-\gamma^2}\right) + \frac{1}{4}\left(\frac{\gamma\zeta + \gamma^4 + \zeta - \gamma^2\zeta}{1-\gamma^2}\right) + \frac{1}{4}\left(\frac{\zeta + \gamma\zeta - \gamma^2\zeta - \gamma^5\zeta + \gamma^6\zeta + \gamma^7}{1-\gamma^2}\right) + \frac{1}{4}V_{S_2}^{T_1} \\
 &= \frac{1}{4}\left(\frac{\zeta + \gamma}{1-\gamma^2}\right) + \frac{1}{4}\left(\frac{\gamma\zeta + \gamma^4 + \zeta - \gamma^2\zeta}{1-\gamma^2}\right) + \frac{1}{4}\left(\frac{\zeta + \gamma\zeta - \gamma^2\zeta - \gamma^5\zeta + \gamma^6\zeta + \gamma^7}{1-\gamma^2}\right) + \frac{1}{4}\left(\frac{\zeta + \gamma^3}{1-\gamma^2}\right) \\
 &= \frac{\zeta(4 + 2\gamma - 2\gamma^2 - \gamma^5 + \gamma^6) + \gamma + \gamma^3 + \gamma^4 + \gamma^7}{4(1-\gamma^2)}
 \end{aligned}$$

**The balanced depth-2 decision tree :** alternates in every state between taking the two available information gathering actions and then a base action. The value of the policy

in the goal state is :

$$\begin{aligned}
 V_G^{\mathcal{T}_2} &= \zeta + \gamma\zeta + \gamma^2 + \gamma^3\zeta + \gamma^4\zeta + \dots \\
 &= \sum_{t=0}^{\infty} \gamma^{3t}\zeta + \sum_{t=0}^{\infty} \gamma^{3t+1}\zeta + \sum_{t=0}^{\infty} \gamma^{3t+2} \\
 &= \frac{\zeta}{1-\gamma^3} + \frac{\gamma\zeta}{1-\gamma^3} + \frac{\gamma^2}{1-\gamma^3}
 \end{aligned}$$

Following the same reasoning for other states we find the objective value for the depth-2 decision tree policy to be :

$$\begin{aligned}
 J(\mathcal{T}_2) &= \frac{1}{4}V_G^{\mathcal{T}_2} + \frac{2}{4}V_{S_2}^{\mathcal{T}_2} + \frac{1}{4}V_{S_1}^{\mathcal{T}_2} \\
 &= \frac{1}{4}V_G^{\mathcal{T}_2} + \frac{2}{4}(\zeta + \gamma\zeta + \gamma^2 0 + \gamma^3 V_G^{\mathcal{T}_2}) + \frac{1}{4}(\zeta + \gamma\zeta + \gamma^2 0 + \gamma^3\zeta + \gamma^4\zeta + \gamma^5 0 + \gamma^6 V_G^{\mathcal{T}_2}) \\
 &= \frac{\zeta(3+3\gamma) + \gamma^2 + \gamma^5 + \gamma^8}{4(1-\gamma^3)}
 \end{aligned}$$

**Infinite tree :** we also consider the infinite tree policy that repeats an information gathering action forever and has objective :  $J(\mathcal{T}_{\text{inf}}) = \frac{\zeta}{1-\gamma}$

**Stochastic policy :** the other non-trivial policy that can be learned by solving a partially observable IBMDP is the stochastic policy that guarantees to reach  $G$  after some time : fifty percent chance to do  $\rightarrow$  and fifty percent chance to do  $\downarrow$ . This stochastic policy has objective value :

$$\begin{aligned}
 V_G^{\text{stoch}} &= \frac{1}{1-\gamma} \\
 V_{S_0}^{\text{stoch}} &= 0 + \frac{1}{2}\gamma V_G^{\text{stoch}} + \frac{1}{2}\gamma V_{S_1}^{\text{stoch}} \\
 V_{S_2}^{\text{stoch}} &= 0 + \frac{1}{2}\gamma V_G^{\text{stoch}} + \frac{1}{2}\gamma V_{S_1}^{\text{stoch}} = V_{S_0}^{\text{stoch}} \\
 V_{S_1}^{\text{stoch}} &= 0 + \frac{1}{2}\gamma V_{S_2}^{\text{stoch}} + \frac{1}{2}\gamma V_G^{\text{stoch}} = \frac{1}{2}\gamma V_{S_0}^{\text{stoch}} + \frac{1}{2}\gamma V_G^{\text{stoch}}
 \end{aligned}$$

Solving these equations :

$$\begin{aligned}
 V_{S_1}^{\text{stoch}} &= \frac{1}{2}\gamma V_{S_0}^{\text{stoch}} + \frac{1}{2}\gamma V_G^{\text{stoch}} \\
 &= \frac{1}{2}\gamma \left( \frac{1}{2}\gamma V_G^{\text{stoch}} + \frac{1}{2}\gamma V_{S_1}^{\text{stoch}} \right) + \frac{1}{2}\gamma V_G^{\text{stoch}} \\
 &= \frac{1}{4}\gamma^2 V_G^{\text{stoch}} + \frac{1}{4}\gamma^2 V_{S_1}^{\text{stoch}} + \frac{1}{2}\gamma V_G^{\text{stoch}} \\
 V_{S_1}^{\text{stoch}} - \frac{1}{4}\gamma^2 V_{S_1}^{\text{stoch}} &= \frac{1}{4}\gamma^2 V_G^{\text{stoch}} + \frac{1}{2}\gamma V_G^{\text{stoch}} \\
 V_{S_1}^{\text{stoch}} \left( 1 - \frac{1}{4}\gamma^2 \right) &= \left( \frac{1}{4}\gamma^2 + \frac{1}{2}\gamma \right) V_G^{\text{stoch}} \\
 V_{S_1}^{\text{stoch}} &= \frac{\frac{1}{4}\gamma^2 + \frac{1}{2}\gamma}{1 - \frac{1}{4}\gamma^2} V_G^{\text{stoch}} \\
 &= \frac{\gamma \left( \frac{1}{4}\gamma + \frac{1}{2} \right)}{1 - \frac{1}{4}\gamma^2} \cdot \frac{1}{1 - \gamma} \\
 &= \frac{\gamma \left( \frac{1}{4}\gamma + \frac{1}{2} \right)}{(1 - \frac{1}{4}\gamma^2)(1 - \gamma)}
 \end{aligned}$$

$$\begin{aligned}
 V_{S_0}^{\text{stoch}} &= \frac{1}{2}\gamma V_G^{\text{stoch}} + \frac{1}{2}\gamma V_{S_1}^{\text{stoch}} \\
 &= \frac{1}{2}\gamma \cdot \frac{1}{1 - \gamma} + \frac{1}{2}\gamma \cdot \frac{\gamma \left( \frac{1}{4}\gamma + \frac{1}{2} \right)}{(1 - \frac{1}{4}\gamma^2)(1 - \gamma)} \\
 &= \frac{\frac{1}{2}\gamma}{1 - \gamma} + \frac{\frac{1}{2}\gamma^2 \left( \frac{1}{4}\gamma + \frac{1}{2} \right)}{(1 - \frac{1}{4}\gamma^2)(1 - \gamma)} \\
 &= \frac{\frac{1}{2}\gamma(1 - \frac{1}{4}\gamma^2) + \frac{1}{2}\gamma^2 \left( \frac{1}{4}\gamma + \frac{1}{2} \right)}{(1 - \frac{1}{4}\gamma^2)(1 - \gamma)} \\
 &= \frac{\frac{1}{2}\gamma - \frac{1}{8}\gamma^3 + \frac{1}{8}\gamma^3 + \frac{1}{4}\gamma^2}{(1 - \frac{1}{4}\gamma^2)(1 - \gamma)} \\
 &= \frac{\frac{1}{2}\gamma + \frac{1}{4}\gamma^2}{(1 - \frac{1}{4}\gamma^2)(1 - \gamma)} \\
 &= \frac{\gamma \left( \frac{1}{2} + \frac{1}{4}\gamma \right)}{(1 - \frac{1}{4}\gamma^2)(1 - \gamma)}
 \end{aligned}$$

$$\begin{aligned}
J(\mathcal{T}_{\text{stoch}}) &= \frac{1}{4}(V_G^{\text{stoch}} + V_{S_0}^{\text{stoch}} + V_{S_1}^{\text{stoch}} + V_{S_2}^{\text{stoch}}) \\
&= \frac{1}{4} \left( \frac{1}{1-\gamma} + 2 \cdot \frac{\gamma(\frac{1}{2} + \frac{1}{4}\gamma)}{(1 - \frac{1}{4}\gamma^2)(1-\gamma)} + \frac{\gamma(\frac{1}{4}\gamma + \frac{1}{2})}{(1 - \frac{1}{4}\gamma^2)(1-\gamma)} \right) \\
&= \frac{1}{4} \left( \frac{1}{1-\gamma} + \frac{2\gamma(\frac{1}{2} + \frac{1}{4}\gamma) + \gamma(\frac{1}{4}\gamma + \frac{1}{2})}{(1 - \frac{1}{4}\gamma^2)(1-\gamma)} \right) \\
&= \frac{1}{4} \left( \frac{1}{1-\gamma} + \frac{\gamma + \frac{1}{2}\gamma^2 + \frac{1}{4}\gamma^2 + \frac{1}{2}\gamma}{(1 - \frac{1}{4}\gamma^2)(1-\gamma)} \right) \\
&= \frac{1}{4} \left( \frac{1}{1-\gamma} + \frac{\frac{3}{2}\gamma + \frac{3}{4}\gamma^2}{(1 - \frac{1}{4}\gamma^2)(1-\gamma)} \right) \\
&= \frac{1}{4} \left( \frac{1 - \frac{1}{4}\gamma^2 + \frac{3}{2}\gamma + \frac{3}{4}\gamma^2}{(1 - \frac{1}{4}\gamma^2)(1-\gamma)} \right) \\
&= \frac{1}{4} \left( \frac{1 + \frac{3}{2}\gamma + \frac{1}{2}\gamma^2}{(1 - \frac{1}{4}\gamma^2)(1-\gamma)} \right) \\
&= \frac{1 + \frac{3}{2}\gamma + \frac{1}{2}\gamma^2}{4(1 - \frac{1}{4}\gamma^2)(1-\gamma)}
\end{aligned}$$

## B.2 Hyperparameters

TABLEAU B.1 – PG Hyperparameter Space (140 combinations)

Hyperparameter	Values	Description
Learning Rate (lr)	0.001, 0.005, 0.01, 0.05, 0.1	Policy gradient step size
Entropy Regularization (tau)	-1.0, -0.1, -0.01, 0.0, 0.01, 0.1, 1.0	Entropy regularization coefficient
Temperature (eps)	0.01, 0.1, 1.0, 10	Softmax temperature
Episodes per Update (n_steps)	20, 200, 2000	Number of episodes per policy update

TABLEAU B.2 – PG-IBMDP Hyperparameter Space (140 combinations)

Hyperparameter	Values	Description
Learning Rate (lr)	0.001, 0.005, 0.01, 0.05, 0.1	Policy gradient step size
Entropy Regularization (tau)	-1.0, -0.1, -0.01, 0.0, 0.01, 0.1, 1.0	Entropy regularization coefficient
Temperature (eps)	0.01, 0.1, 1.0, 10	Softmax temperature
Episodes per Update (n_steps)	10, 100, 1000	Number of episodes per policy update

TABLEAU B.3 – QL Hyperparameter Space (192 combinations)

Hyperparameter	Values	Description
Epsilon Schedules	(0.3, 1), (0.3, 0.99), (1, 1)	Initial exploration and decrease rate
Epsilon Schedules	(0.1, 1), (0.1, 0.99), (0.3, 0.99)	Initial exploration and decrease rate
Lambda	0.0, 0.3, 0.6, 0.9	Eligibility trace decay
Learning Rate (lr_o)	0.001, 0.005, 0.01, 0.1	Observation Q-learning rate
Optimistic	True, False	Optimistic initialization

TABLEAU B.4 – QL-Asym Hyperparameter Space (768 combinations)

Hyperparameter	Values	Description
Epsilon Schedules	(0.3, 1), (0.3, 0.99), (1, 1)	Initial exploration and decrease rate
Epsilon Schedules	(0.1, 1), (0.1, 0.99), (0.3, 0.99)	Initial exploration and decrease rate
Lambda	0.0, 0.3, 0.6, 0.9	Eligibility trace decay
Learning Rate (lr_o)	0.001, 0.005, 0.01, 0.1	Observation Q-learning rate
Learning Rate (lr_v)	0.001, 0.005, 0.01, 0.1	State-action Q-learning rate
Optimistic	True, False	Optimistic initialization

TABLEAU B.5 – QL-IBMDP Hyperparameter Space (192 combinations)

Hyperparameter	Values	Description
Epsilon Schedules	(0.3, 1), (0.3, 0.99), (1, 1)	Initial exploration and decrease rate
Epsilon Schedules	(0.1, 1), (0.1, 0.99), (0.3, 0.99)	Initial exploration and decrease rate
Lambda	0.0, 0.3, 0.6, 0.9	Eligibility trace decay
Learning Rate (lr_v)	0.001, 0.005, 0.01, 0.1	State-action Q-learning rate
Optimistic	True, False	Optimistic initialization

TABLEAU B.6 – SARSA Hyperparameter Space (192 combinations)

Hyperparameter	Values	Description
Epsilon Schedules	(0.3, 1), (0.3, 0.99), (1, 1)	Initial exploration and decrease rate
Epsilon Schedules	(0.1, 1), (0.1, 0.99), (0.3, 0.99)	Initial exploration and decrease rate
Lambda	0.0, 0.3, 0.6, 0.9	Eligibility trace decay
Learning Rate (lr_o)	0.001, 0.005, 0.01, 0.1	Observation SARSA learning rate
Optimistic	True, False	Optimistic initialization

TABLEAU B.7 – SARSA-Asym Hyperparameter Space (768 combinations)

Hyperparameter	Values	Description
Epsilon Schedules	(0.3, 1), (0.3, 0.99), (1, 1)	Initial exploration and decrease rate
Epsilon Schedules	(0.1, 1), (0.1, 0.99), (0.3, 0.99)	Initial exploration and decrease rate
Lambda	0.0, 0.3, 0.6, 0.9	Eligibility trace decay
Learning Rate (lr_o)	0.001, 0.005, 0.01, 0.1	Observation SARSA learning rate
Learning Rate (lr_v)	0.001, 0.005, 0.01, 0.1	State-action SARSA learning rate
Optimistic	True, False	Optimistic initialization

TABLEAU B.8 – SARSA-IBMDP Hyperparameter Space (192 combinations)

Hyperparameter	Values	Description
Epsilon Schedules	(0.3, 1), (0.3, 0.99), (1, 1)	Initial exploration and decrease rate
Epsilon Schedules	(0.1, 1), (0.1, 0.99), (0.3, 0.99)	Initial exploration and decrease rate
Lambda	0.0, 0.3, 0.6, 0.9	Eligibility trace decay
Learning Rate (lr_v)	0.001, 0.005, 0.01, 0.1	State-action SARSA learning rate
Optimistic	True, False	Optimistic initialization





# Table des matières

<b>Résumé</b>	<b>vii</b>
<b>Sommaire</b>	<b>ix</b>
<b>Preliminary Concepts</b>	<b>1</b>
Interpretable Sequential Decision Making . . . . .	1
What is Sequential Decision Making? . . . . .	1
What is Interpretability? . . . . .	2
What are existing approaches for learning interpretable programs? .	5
Other related works . . . . .	7
Explainability . . . . .	7
Misalignement . . . . .	7
Programmatic RL . . . . .	7
Causality . . . . .	7
Mechanistic Interpretability in the era of large language models . . .	7
Outline of the Thesis . . . . .	7
Technical Preliminaries . . . . .	8
What are decision trees? . . . . .	8
How to learn decision trees? . . . . .	9
Markov decision processes/problems . . . . .	11
Exact solutions for Markov decision problems . . . . .	12
Reinforcement learning of approximate solutions to MDPs . . . . .	13
Imitation learning, the baseline (indirect) global interpretable reinforcement learning method . . . . .	16
<b>I A Difficult Problem : Direct Interpretable Reinforcement Learning</b>	<b>17</b>
<b>1 A Formal Framework for the Reinforcement Learning of Decision Tree Policies</b>	<b>19</b>
1.1 Learning Decision Tree policies . . . . .	19
1.2 Example of the limitations of indirect methods . . . . .	20
1.3 Iterative Bounding Markov Decision Processes . . . . .	20

1.3.1	Formalism . . . . .	22
1.3.2	From Policies to Trees . . . . .	23
1.3.3	Didactic example . . . . .	24
1.4	The promise of IBMDPs . . . . .	25
<b>2</b>	<b>Limits of Direct (Deep) Reinforcement Learning of Decision Tree Policies</b>	<b>29</b>
2.1	Reproducing “Iterative Bounding MDPs : Learning Interpretable Policies via Non-Interpretable Methods” . . . . .	29
2.1.1	IBMDP formulation . . . . .	30
2.1.2	Modified Deep Reinforcement Learning algorithms . . . . .	30
2.2	Experimental setup . . . . .	31
2.2.1	(IB)MDP . . . . .	31
2.2.2	Baselines . . . . .	32
2.2.3	Metrics . . . . .	34
2.3	Results . . . . .	35
2.3.1	Can Modified Deep Reinforcement Learning consistently learn good decision tree policies for CartPole? . . . . .	35
2.3.2	How does direct interpretable reinforcement learning perform compared to the indirect approach? . . . . .	38
2.4	Discussion . . . . .	40
<b>3</b>	<b>Understanding the Limits of Direct Reinforcement Learning</b>	<b>41</b>
3.0.1	Partially Observable IBMDPs . . . . .	41
3.1	Constructing POIBMDPs which optimal solutions are the depth-1 tree . . . . .	44
3.2	Results . . . . .	48
3.2.1	Experimental Setup . . . . .	48
3.2.2	Can RL baselines solve POIBMDPs? . . . . .	49
3.2.3	How difficult is it to learn in POIBMPDs? . . . . .	52
3.3	Conclusion . . . . .	55
<b>4</b>	<b>When transitions are uniform POIBMDPs are fully observable</b>	<b>57</b>
4.1	How well can RL baselines learn in Classification POIBMDPs? . . . . .	59
<b>II</b>	<b>An easier problem : Learning Decision Trees for MDPs that are Classification tasks</b>	<b>61</b>
<b>5</b>	<b>DPDT-intro</b>	<b>63</b>
5.1	Introduction . . . . .	64
5.2	Related Work . . . . .	65
<b>6</b>	<b>DPDT-paper</b>	<b>67</b>
6.1	Decision Trees for Supervised Learning . . . . .	67
6.2	Decision Tree Induction as an MDP . . . . .	68

6.3	Algorithm . . . . .	69
6.3.1	Constructing the MDP . . . . .	69
6.3.2	Heuristic splits generating functions . . . . .	69
6.3.3	Dynamic Programming to solve the MDP . . . . .	71
6.3.4	Performance Guarantees of DPDT . . . . .	72
6.3.5	Proof of Improvement over CART . . . . .	73
6.3.6	Practical Implementation . . . . .	74
6.4	Empirical Evaluation . . . . .	75
6.4.1	DPDT optimizing capabilities . . . . .	75
6.4.2	DPDT generalization capabilities . . . . .	80
6.5	Application of DPDT to Boosting . . . . .	83
6.5.1	Boosted-DPDT . . . . .	83
6.5.2	(X)GB-DPDT . . . . .	83
6.6	Proof of Proposition 2 . . . . .	85
6.7	Additional Experiments and Hyperparameters . . . . .	85
<b>7</b>	<b>Conclusion</b>	<b>89</b>
7.1	Conclusion . . . . .	89
7.2	What about imitation? . . . . .	89
<b>III</b>	<b>Beyond Decision Trees : what can be done with other Interpretable Policies?</b>	<b>91</b>
<b>8</b>	<b>Imitation and Evaluation</b>	<b>93</b>
8.1	Intro . . . . .	93
<b>9</b>	<b>Evaluation</b>	<b>95</b>
9.1	Methodology Overview . . . . .	95
9.2	Computing Baseline Policies . . . . .	97
9.2.1	Setup . . . . .	97
9.2.2	Ablation study of imitation learning . . . . .	101
9.3	Measuring Policy Interpretability . . . . .	102
9.3.1	From Policy to Program . . . . .	102
9.3.2	Interpretability-performance trade-offs . . . . .	104
9.3.3	Verifying interpretable policies . . . . .	105
9.4	Experimental details . . . . .	107
9.5	All interpretability-performance trade-offs . . . . .	107
<b>10</b>	<b>Conclusion Imitation</b>	<b>111</b>
10.1	Limitations and conclusions . . . . .	111
	<b>Conclusion générale</b>	<b>113</b>
	<b>Bibliographie</b>	<b>115</b>

<b>A Programmes informatiques</b>	<b>125</b>
<b>B Appendix I</b>	<b>127</b>
B.1 Tree value computations . . . . .	127
B.2 Hyperparameters . . . . .	131
<b>Table des matières</b>	<b>135</b>