

Limits of reinforcement learning for decision trees in Markov decision processes

Anonymous authors

Paper under double-blind review

Keywords: Reinforcement learning, decision trees, POMDP

Summary

The summary appears on the cover page. Although it can be identical to the abstract, it does not have to be. One might choose to omit the stated contributions in the Summary, given that they will be stated in the box below. The original abstract may also be extended to two paragraphs. The authors should ensure that the contents of the cover page fit entirely on a single page. The cover page does **not** count towards the 8–12 page limit.

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.

Contribution(s)

1. Provide a succinct but precise list of the contribution(s) of the paper. Use contextual notes to avoid implications of contributions more significant than intended and to clarify and situate the contribution relative to prior work (see the examples below). If there is no additional context, enter “None”. Try to keep each contribution to a single sentence, although multiple sentences are allowed when necessary. If using complete sentences, include punctuation. If using a single sentence fragment, you may omit the concluding period. A single contribution can be sufficient, and there is no limit on the number of contributions. Submissions will be judged mostly on the contributions claimed on their cover pages and the evidence provided to support them. Major contributions should not be claimed in the main text if they do not appear on the cover page. Overclaiming can lead to a submission being rejected, so it is important to have well-scoped contribution statements on the cover page.

Context: None

2. The submission template for submissions to RLJ/RLC 2026

Context: Built from previous RLC/RLJ, ICLR, and TMLR submission templates

3. *[Example of one contribution and corresponding contextual note for the paper “Policy gradient methods for reinforcement learning with function approximation” (?).]*

This paper presents an expression for the policy gradient when using function approximation to represent the action-value function.

Context: Prior work established expressions for the policy gradient without function approximation (?).

Limits of reinforcement learning for decision trees in Markov decision processes

Anonymous authors

Paper under double-blind review

Abstract

For applications like medicine, machine learning models ought to be interpretable. In that case, models like decision trees are preferred over neural networks because humans can read their predictions from the root to the leaves. Learning such decision trees for sequential decision making problems is a relatively new research direction and most of the existing literature focuses on imitating (or distilling) neural networks. In contrast, we study reinforcement learning (RL) algorithms that *directly* return decision trees optimizing some trade-off of cumulative rewards and interpretability in a Markov decision process (MDP). We show that such algorithms can be seen as learning policies for partially observable Markov decision processes (POMDPs). We use this parallel to understand why in practice it is often easier to use imitation learning than to learn the decision tree from scratch for MDPs.

1 Introduction

Interpretability in machine learning is commonly divided into local and global approaches [Glanois et al. \(2024\)](#). Local methods—also referred to as explainability or post-hoc methods [Lipton \(2018\)](#)—provide explanations for individual predictions using tools such as local linear approximations [Ribeiro et al. \(2016\)](#), saliency maps [Puri et al. \(2020\)](#), feature attributions [Lundberg & Lee \(2017\)](#), or attention mechanisms [Shi et al. \(2022\)](#). Although widely used, these methods approximate the behavior of an underlying black-box model and may therefore be unfaithful to its true computations [Atrey et al. \(2020\)](#).

Global interpretability approaches instead restrict the model class so that the learned model is transparent by construction. Decision trees [Breiman et al. \(1984\)](#) are a canonical example, as their predictions can be inspected, reasoned about, and formally verified. This makes them particularly attractive for safety-critical applications and has motivated extensive research in supervised learning [Murthy & Salzberg \(1995\)](#); [Verwer & Zhang \(2019\)](#); [Demirovic et al. \(2022\)](#); [Demirović et al. \(2023\)](#); [van der Linden et al. \(2023\)](#).

Extending global interpretability to sequential decision making, however, remains challenging. Existing approaches largely rely on *indirect* methods [Milani et al. \(2024\)](#): a high-performing but opaque policy (typically a neural network) is first learned using reinforcement learning, and an interpretable model is then trained to imitate its behavior. A prominent example is VIPER [Bastani et al. \(2018\)](#), which distills neural network policies into decision trees using imitation learning [Ross et al. \(2010\)](#). Such methods have demonstrated strong empirical performance and enable formal verification [Wu et al. \(2024\)](#), but they optimize a surrogate objective—policy imitation—rather than the original reinforcement learning objective. As a result, the best decision tree policy for the task may differ substantially from the tree that best approximates a neural expert. The curious reader will find an example of this phenomenon in the appendix 8.

This limitation motivates the study of *direct* approaches that learn interpretable policies by optimizing the reinforcement learning objective itself. While direct decision tree learning is well understood

in supervised settings, it is far less developed for sequential decision making. Understanding why direct optimization is difficult—and when it can succeed—is the central focus of this work.

In this article, we show that reinforcement learning of decision tree policies for MDPs, i.e. learning a decision tree that directly optimizes the cumulative reward of the process without relying on a black-box expert, is often very difficult. To do so, we construct very simple MDPs for which we know optimal decision tree policies and show that RL consistently fails to retrieve those policies. We identify partial observability as a key reason for those failures.

In section 2, we present the related work on reinforcement learning to train decision tree policies for MDPs. In section 3, we present key concepts for decision trees, MDPs, and the formalism of [Topin et al. \(2021\)](#) for reinforcement learning of decision tree policies. In section 4, we show that this direct approach is equivalent to learning a *deterministic memoryless* policy for partially observable MDP (POMDP) [Sondik \(1978\)](#) which is a hard problem [Littman \(1994\)](#). In section 5, we present our methodology to benchmark RL algorithms that train decision tree policies. In section 5.3, we show that when RL fails to retrieve optimal decision tree policies for MDPs it is most likely because partial observability is involved.

2 Related work

There exist reinforcement learning algorithms that directly train decision tree policies optimizing the cumulative rewards in a given MDP. These approaches can be divided into methods based on *parametric* and *non-parametric* trees.

Parametric decision trees fix the tree structure *a priori*—including depth, node arrangement, and selected state features—and only learn the decision thresholds. This formulation enables differentiability and allows direct optimization of the RL objective using policy gradient methods [Sutton et al. \(1999\)](#). Several works [Silva et al. \(2020\)](#); [Vos & Verwer \(2024\)](#); [Marton et al. \(2025\)](#) employ PPO to train such differentiable trees. While these methods can achieve strong performance, they require the tree structure to be specified in advance, making it difficult to adaptively trade off interpretability and performance. An overly complex structure may require post-hoc pruning, whereas an insufficiently expressive structure may fail to represent good policies. Moreover, [Marton et al. \(2025\)](#) reports that additional stabilization techniques, such as adaptive batch sizes, are often necessary for direct reinforcement learning to match indirect imitation methods [Bastani et al. \(2018\)](#) performances. Non-parametric decision trees, by contrast, are the standard model in supervised learning, where algorithms efficiently construct trees that balance predictive performance and interpretability [Breiman et al. \(1984\)](#); [Bertsimas & Dunn \(2017\)](#). This trade-off is optimized by using a regularized training objective. However, training non-parametric trees to trade-off MDP cumulative rewards and interpretability is largely unexplored [Milani et al. \(2024\)](#). To the best of our knowledge, the only work that studies this setting: [Topin et al. \(2021\)](#). Topin et al. introduce *iterative bounding MDPs* (IBMDPs), which augment the downstream MDP with additional state features, actions, rewards, and transitions. They show that certain policies in the IBMDP correspond to decision tree policies for the downstream MDP. Hence, standard RL algorithms can be used to learn such policies in IBMDPs. Finally, a few specialized methods exist for restricted problem classes. For maze-like MDPs, [Mansour et al. \(2022\)](#) proves the existence of optimal decision tree policies and provides a constructive algorithm. In settings where the MDP model is fully known, [Vos & Verwer \(2023\)](#) use planning to compute shallow parametric decision tree policies. Next, we recall useful technical material.

3 Technical preliminaries

3.1 Markov decision processes

Markov decision processes were first introduced in the 1950s by Richard Bellman [Bellman \(1957\)](#). Informally, an MDP models how an agent acts over time to achieve a goal. At every time step,

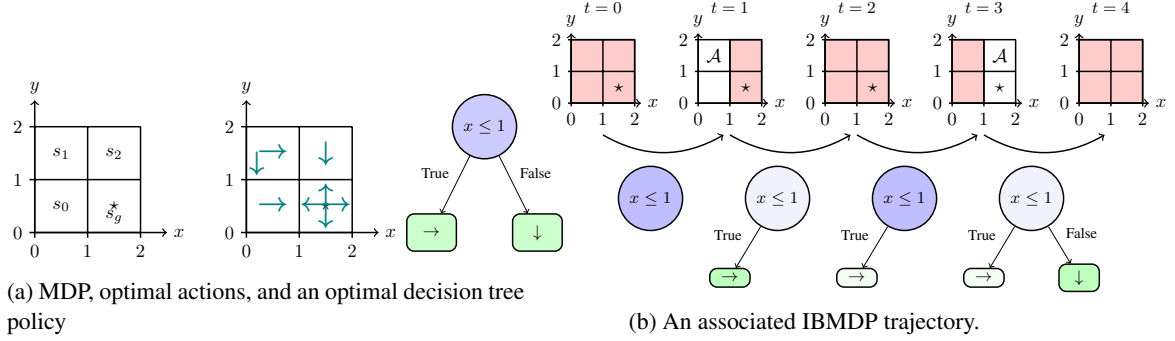


Figure 1: On the left of figure 1a, a grid world MDP with, four states, four directional actions that moves an agent, and rewards of 0 for every transitions except when taking an action in the bottom right goal state (\star). The states have discrete labels that represent their coordinates in the $[0, 2] \times [0, 2]$ square, e.g. for s_0 , $x = 0.5$, $y = 0.5$. In the centre, an optimal actions w.r.t the RL objective (definition 2). On the right, an optimal depth-1 decision tree policy w.r.t the RL objective (definition 2) that takes a state label as input, performs some tests on the state coordinate, and returns a directional action. On figure 1b, an IBMDP trajectory when the downstream MDP is the grid world from figure 1a. In the top row, we graphically represent what the IBMDP state at a given time t , and in the bottom row, we present the corresponding decision tree policy traversal. When the pink covers the whole grid, the information contained in the observation \mathbf{o}_t could be interpreted as ‘the current state features could be anywhere in the grid’. The more information-gathering actions are taken, the more refined the bounds on the current downstream state features get. At $t = 0$, the downstream state features are $s_0 = (0.5, 1.5)$ and the initial observation is always the downstream MDP default state feature bounds (definition 5): $\mathbf{o}_0 = (0, 2, 0, 2)$ because the downstream state features are in $[0, 2] \times [0, 2]$. This means that the overall IBMDP state is $\mathbf{s}_{IB} = (0.5, 1.5, 0, 2, 0, 2)$. The first action is an IGA $\langle x, 1 \rangle$ that tests the feature x of the downstream state against the value 1 and the reward is ζ (definition 5). This transition corresponds to going through an internal node $x \leq 1$ in a decision tree policy as illustrated in the figure. At $t = 1$, after gathering the information that the x -value of the current downstream state is below 1, the observation is updated with the refined bounds $\mathbf{o}_1 = (0, 1, 0, 2)$, i.e. more information has been gathered and the obstructed pink area shrinks. The downstream state features remain unchanged. The agent then takes a downstream action that is to move right. This gives a reward 0, resets the observation to the original downstream state feature bounds, and changes the features to $s_2 = (1.5, 1.5)$. And the trajectory continues like this until the absorbing downstream state $s_4 = (1.5, 0.5)$ is reached. By masking the current downstream state features to an agent, we would force it to take information-gathering actions otherwise it would not know how to act optimally (figure 1a).

the agent observes its current state (e.g., patient weight and tumor size) and takes an action (e.g., administers a certain amount of chemotherapy). The agent receives a reward that helps evaluate the quality of the action with respect to the goal (e.g., tumor size decrease when the objective is to cure cancer). Finally, the agent transitions to a new state (e.g., the updated patient state) and repeats this process over time:

Definition 1 (Markov decision process). An MDP is a tuple $\mathcal{M} = \langle S, A, R, T, T_0 \rangle$. S is a finite set of states representing all possible configurations of the environment. A is a finite set of actions available to the agent. $R : S \times A \rightarrow \mathbb{R}$ is a deterministic reward function that assigns a real-valued reward to each state-action pair. While in general reward functions are often stochastic, in this manuscript we focus deterministic ones without loss of generality. $T : S \times A \rightarrow \Delta(S)$ is the transition function that maps state-action pairs to probability distributions over next states $\Delta(S)$. $T_0 \in \Delta(S)$ is the initial distribution over states.

Informally, we would like to act in an MDP so that we obtain as much reward as possible over time. We can formally define this objective, that we call the reinforcement learning objective, as follows:

Definition 2 (Reinforcement learning objective). Given an MDP (definition 1) $\mathcal{M} \equiv \langle S, A, R, T, T_0 \rangle$, the goal of reinforcement learning for sequential decision making is to find a model, also known as a policy, $\pi : S \rightarrow A$ that maximizes the expected discounted sum of rewards:

$$J(\pi) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t R(s_t, \pi(s_t)) \mid s_0 \sim T_0, s_{t+1} \sim T(s_t, \pi(s_t)) \right]$$

where $0 < \gamma \leq 1$ is the discount factor that controls the trade-off between immediate and future rewards.

Algorithms presented in this article aim to find an optimal policy $\pi^* \in \operatorname{argmax}_{\pi} J(\pi)$ that maximizes the above reinforcement learning objective. In particular, RL algorithms [Sutton & Barto \(1998\)](#); [Sutton et al. \(1999\)](#); [Watkins & Dayan \(1992\)](#); [Mnih et al. \(2015\)](#); [Schulman et al. \(2017\)](#) learn such optimal policies using data of MDP interactions without prior knowledge of the reward and transition models. Useful quantities for such algorithms include *value* of states and actions.

Definition 3 (Value of a state). In an MDP \mathcal{M} (definition 1), the value of a state $s \in S$ under policy π is the expected discounted sum of rewards starting from state s and following policy π :

$$V^{\pi}(s) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t R(s_t, \pi(s_t)) \mid s_0 = s, s_{t+1} \sim T(s_t, \pi(s_t)) \right]$$

Applying the Markov property gives a recursive definition of the value of s under policy π : $V^{\pi}(s) = R(s, \pi(s)) + \gamma \mathbb{E}[V^{\pi}(s') \mid s' \sim T(s, \pi(s))]$. The optimal value of a state $s \in S$, $V^*(s)$, is the value of state s when following the optimal policy π^* (the policy that maximizes the RL objective (definition 2)): $V^*(s) = V^{\pi^*}(s)$. Similarly, the optimal value of a state-action pair $(s, a) \in S \times A$, $Q^*(s, a)$, is the value when taking action a in state s and then following the optimal policy: $Q^*(s, a) = R(s, a) + \gamma \mathbb{E}[V^*(s') \mid s' \sim T(s, a)]$.

3.2 Decision tree policies

While other interpretable policy classes exist [Kohler et al. \(2025b\)](#), one conjecture from [Glanois et al. \(2024\)](#) is that interpretable models are all hard to optimize or learn because they are non-differentiable in nature. This is something that will be key in our study of decision tree policy that we introduce next.

Definition 4 (Decision tree policy). A decision tree policy is a rooted tree $\pi_{\mathcal{T}} = (\mathcal{N}, E)$. Each internal node $\nu \in \mathcal{N}$ is associated with a test that maps an MDP state attribute to a Boolean, e.g. $s_i \leq v$ for $v \in \mathbb{R}$. Each edge $e \in E$ from an internal node corresponds to an outcome of the associated test function. Each leaf node $l \in \mathcal{N}$ is associated with an MDP action $a_l \in A$. For any input $s \in S$, the tree defines a unique path from root to leaf, determining the prediction $\pi_{\mathcal{T}}(s) = a_l$ where l is the reached leaf. The depth of a tree is the maximum path length from root to any leaf.

In figure 1a, we present an example Markov decision process (definition 1), the optimal actions that maximize the RL objective (definition 2), and a decision tree policy (definition 4) that also maximizes the RL objective. Next, we present the class of MDPs introduced in Topin et al. (2021) useful for our to understand reinforcement learning of decision tree policies that directly optimize the RL objective in an MDP.

3.3 Iterative bounding Markov decision processes

The key thing to know about iterative bounding Markov decision processes (IBMDPs) is that they are, as their name suggests, MDPs (definition 1). Hence, IBMDPs admit an optimal deterministic Markovian policy that maximizes the RL objective Bellman (1957). From now on, we will assume that all the MDPs we consider are MDPs with continuous state spaces and a finite set of actions, and we use bold fonts for states and observations as that are vector-valued. However all our results generalize to discrete states (in \mathbb{Z}^m) MDPs that we can factor using one-hot encodings. Given an MDP for which we want to learn a decision tree policy—the downstream MDP—IBMDP states are concatenations of the downstream MDP state features and some observations. Those observations are information about the downstream state features that are refined—‘iteratively bounded’—at each step. Those observations essentially represent some knowledge about where some downstream state features lie in the state space. Actions available in an IBMDP are: (i) the actions of the downstream MDP, that change downstream state features, and (ii) *information-gathering* actions (IGAs) that change the aforementioned observations. Now, downstream actions in an IBMDP are rewarded like in the downstream MDP, this ensures that the RL objective w.r.t. the downstream MDP is encoded in the IBMDP reward. When taking an information-gathering action, the reward is an arbitrary value such that optimizing the RL objective in the IBMDP is equivalent to optimizing some trade-off between interpretability and the RL objective in the downstream MDP. Before showing how to get decision tree policies from IBMDP policies, we give a formal definition of IBMDPs following Topin et al. (2021):

Definition 5 (Iterative bounding Markov decision process (Topin et al. (2021), section 4.1)). Given a downstream MDP $\mathcal{M} \equiv \langle S, A, R, T, T_0 \rangle$ (definition 1), an associated iterative bounding Markov decision process \mathcal{M}_{IB} is a tuple:

$$\langle \overbrace{S \times O}^{\text{State space}}, \overbrace{A \cup A_{info}}^{\text{Action space}}, \overbrace{(R, \zeta)}^{\text{Reward}}, \overbrace{(T_{info}, T, T_0)}^{\text{Transitions}} \rangle$$

S are the downstream MDP state features. Downstream state features $\mathbf{s} = (s_1, \dots, s_p) \in S$ are bounded: $s_j \in [L_j, U_j]$ where $-\infty < L_j \leq U_j < \infty \forall 1 \leq j \leq p$. O are observations. They represent bounds on the downstream state features: $O \subsetneq S^2 = [L_1, U_1] \times \dots \times [L_p, U_p] \times [L_1, U_1] \times \dots \times [L_p, U_p]$. So the complete IBMDP state space is $S \times O$: the concatenations of downstream state features and observations. Given some downstream state features $\mathbf{s} = (s_1, \dots, s_p) \in S$ and some observation $\mathbf{o} = (L_1, U_1, \dots, L_p, U_p)$, an IBMDP state

is $\mathbf{s}_{IB} = (\overbrace{s_1, \dots, s_p}^{\text{downstream state features}}, \overbrace{L_1, U_1, \dots, L_p, U_p}^{\text{observation}})$. A are the downstream MDP actions. A_{info} are *information-gathering* actions of the form $\langle j, v \rangle$ where j is a state feature index $1 \leq j \leq p$ and v is a real number between L_j and U_j . So the complete action space of an IBMDP is the set of downstream MDP actions and information-gathering actions $A \cup A_{info}$. $R : S \times A \rightarrow \mathbb{R}$ is the downstream MDP reward function. ζ is a reward signal for taking an information-gathering action. So the IBMDP reward function is to get a reward from the downstream MDP if the action is a downstream MDP action or to get ζ if the action is an IGA. $T_{info} : S \times O \times (A_{info} \cup A) \rightarrow \Delta(S \times O)$ is the transition function of IBMDPs: given some observation $\mathbf{o}_t = (L'_1, U'_1, \dots, L'_p, U'_p) \in O$ and downstream state features $\mathbf{s}_t = (s'_1, s'_2, \dots, s'_p)$ if an IGA $\langle j, v \rangle$ is taken, the new observation \mathbf{o}_{t+1} is $(L'_1, U'_1, \dots, L'_j, \min\{v, U'_j\}, \dots, L'_p, U'_p)$ if $s_j \leq v$ or $(L'_1, U'_1, \dots, \max\{v, L'_j\}, U'_j, \dots, L'_p, U'_p)$ if $s_j > v$. If a downstream action is taken, the observation is reset to the default downstream state feature bounds $(L_1, U_1, \dots, L_p, U_p)$ and the downstream state features change according to the downstream MDP transition function: $\mathbf{s}_{t+1} \sim T(\mathbf{s}_t, a_t)$.

At initialization, the downstream state features are drawn from the downstream MDP initial distribution T_0 and the observation is always set to the default downstream state features bounds $\mathbf{o}_0 = (L_1, U_1, \dots, L_p, U_p)$.

We present an IBMDP for a the grid-world MDP (figure 1a) in figure 1b. Now remains to extract a decision tree policy (definition 4) for a downstream MDP \mathcal{M} from a policy for an associated IBMDP \mathcal{M}_{IB} .

3.4 From policies to trees

information-gathering actions (definition 5) resemble the tests $1_{\{x_{-,j} \leq v\}}$ that make up internal decision tree nodes (figure 1b). Indeed, a policy taking actions in an IBMDP essentially builds a tree by taking sequences of IGAs (internal nodes) and then a downstream action (leaf node) and repeats this process over time. In particular, the IGA rewards ζ can be seen as a regularization or a penalty for interpretability: if ζ is very small compared to downstream rewards, a policy will try to take downstream actions as often as possible, i.e. build shallow trees with short paths between root and leaves.

Topin et al. (2021) show that not all IBMDP policies are decision tree policies for the downstream MDP. In particular, their algorithm that converts IBMDP policies into decision trees (algorithm 1) takes as input deterministic policies depending solely on the observations of the IBMDP. If the policies were not deterministic, different subtrees could stem from similar IBMDP observations: this means that the corresponding policy would be a stochastic decision tree Blanc et al. (2021). While there is nothing wrong in learning stochastic decision tree policies from a performance standpoint, interpreting a stochastic policy is an open problem that is not our focus. If the policies were depending on the current full IBMDP state rather than solely on the current IBMDP observation, then a learning agent has 0 incentive to take information-gathering actions to build a decision tree policy as all the state information required to optimally control the downstream MDP as well as downstream actions are available to the agent (definition 5 and figure 1b). The connections between partially observable MDPs (POMDPs Sondik (1978); Sigaud & Buffet (2013)) and extracting decision tree policies from IBMDPs might seem obvious but they are absent from the original IBMDP paper Topin et al. (2021) and from subsequent work. In the next section we bridge this gap.

4 Bridging the gap with the partially observable MDPs literature

4.1 An adequate formalism

To better understand reinforcement learning of decision tree policies for MDPs, we explicitly rewrite the problem of optimizing a deterministic policy depending on current observations in an IBMDP as the problem of optimizing a deterministic policy depending only on current observations—also known as a deterministic *memoryless* policy—in a partially observable Markov decision process (POMDP Sondik (1978)). By doing so, we can leverage results from the POMDP literature that is richer than interpretable reinforcement learning literature.

Definition 6 (Partially observable Markov decision process). A partially observable Markov decision process is a tuple $\langle X, A, O, R, T, T_0, \Omega \rangle$. X is the hidden state space. A is a finite set of actions. O is a set of observations. $T : X \times A \rightarrow \Delta(X)$ is the transition function, where $T(\mathbf{x}_t, a, \mathbf{x}_{t+1}) = P(\mathbf{x}_t | \mathbf{x}_{t+1}, a)$ is the probability of transitioning to state \mathbf{x}_t when taking action a in state \mathbf{x} . T_0 is the initial distribution over states. $\Omega : X \rightarrow \Delta(O)$ is the observation function, where $\Omega(\mathbf{o}, a, \mathbf{x}) = P(\mathbf{o} | \mathbf{x}, a)$ is the probability of observing \mathbf{o} in state \mathbf{x} . $R : X \times A \rightarrow \mathbb{R}$ is the reward function, where $R(\mathbf{x}, a)$ is the immediate reward for taking action a in state \mathbf{x} . Note that $\langle X, A, R, T, T_0 \rangle$ defines an MDP (definition 1).

Next, we can define partially observable iterative bounding Markov decision processes (POIBMDPs). They are IBMDPs (definition 5) for which we explicitly define an observation space and an observation function.

Definition 7 (Partially observable iterative bounding Markov decision process). a partially observable iterative bounding Markov decision process \mathcal{M}_{POIB} is a tuple:

$$\langle \overbrace{S \times O}^{\text{States}}, \overbrace{A \cup A_{info}}^{\text{Action space}}, \overbrace{O}^{\text{Observations}}, \overbrace{(R, \zeta)}^{\text{Rewards}}, \overbrace{(T_{info}, T, T_0), \Omega}^{\text{Transitions}} \rangle$$

, where $\langle S \times O, A \cup A_{info}, (R, \zeta), (T, T_0, T_{info}) \rangle$ is an IBMDP (definition 5). The transition function Ω maps concatenation of state features and observations—IBMDP states—to observations, $\Omega : S \times O \rightarrow O$, with $P(o|(s, o)) = 1$

POIBMDPs are particular instances of POMDPs where the observation function simply applies a mask over some features of the hidden state. This setting has other names in the literature. For example, POIBMDPs are mixed observability MDPs Araya-López et al. (2010) with downstream MDP state features as the *hidden variables* and feature bounds as *visible* variables. POIBMDPs can also be seen as non-stationary MDPs (N-MDPs) Singh et al. (1994) in which there is one different transition function per downstream MDP state: these are called hidden-mode MDPs Choi et al. (2001). Following Singh et al. (1994) we can write the value of a deterministic memoryless policy $\pi : O \rightarrow A \cup A_{info}$ in observation o .

Definition 8 (Value of an observation). In a POIBMDP (definition 7), the expected cumulative discounted reward of a deterministic memoryless 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 section 4 from Singh et al. (1994) for the full definition) of the hidden POIBMDP state (s, o') given the partial observation o and $V^\pi((s, o'))$ the classical state-value function (definition 3). We abuse notation and denote both values of observations and values of states by V since the function input is not ambiguous.

4.2 Reinforcement learning in POMDP

In general, the policy that maximizes the RL objective (definition 2) in a POMDP (definition 6) maps “belief states” or observation histories to actions Sigaud & Buffet (2013). Hence, those policies do not correspond to decision trees since we require that policies depend only on the current observation (algorithm 1). If we did not have this constraint, we could apply any standard RL algorithm to solve POIBMDPs by seeking policies depending on belief states or observations histories because those are sufficient to optimally control any POMDP Sigaud & Buffet (2013); Lambrechts et al. (2025a).

In particular, the problem of finding the optimal deterministic memoryless policies for POMDPs is NP-HARD, even with full knowledge of transitions and rewards (section 3.2 from Littman (1994)). It means that it is impractical to enumerate all possible policies and take the best one. For even moderate-sized POMDPs, a brute-force approach would take a very long time since there are $|A|^{|O|}$ deterministic memoryless policies. Hence it is interesting to study reinforcement learning for finding the best deterministic memoryless policy since it would not enumerate the whole solution space.

In Singh et al. (1994), the authors show that the optimal memoryless policy can be stochastic. Hence, policy gradient algorithms Sutton et al. (1999)—that return stochastic policies—are to avoid since we seek the best *deterministic* policy. Furthermore, the optimal deterministic memoryless policy might not maximize all the values of all observations simultaneously Singh et al. (1994) which makes it difficult to use TD-learning algorithms like Q-learning Watkins & Dayan (1992). Indeed, doing a TD-learning update of one observation’s value (definition 8) can change the value of *all* other observations in an uncontrollable manner because of the dependence in $P^\pi((s, o')|o)$ (definition 8).

Despite those hardness results, applying RL to POMDPs, by naively replacing the processes (hidden) states x by the observation o in Q-learning or Sarsa Sutton & Barto (1998), has already demonstrated successful in practice Loch & Singh (1998). More recently, the framework Baisero et al. (2022); Baisero & Amato (2022) called asymmetric RL, has also shown promising results to learn policies for POMDPs. Asymmetric RL algorithms train a model—a policy or a value function—depending

on hidden state (only available at train time) and a history dependent (or observation dependent) model. The history or observation dependent model serves as target or critic to train the hidden state dependent model. The history dependent (or observation dependent) model can thus be deployed in the POMDP after training since it does not require access to the hidden state to output actions. In appendix 12 we present asymmetric variants of standard tabular RL algorithm. For example, asymmetric Q-learning is a variant of Q-learning that returns a deterministic memoryless policy. Given a POMDP, asymmetric Q-learning trains a partially observable Q-function $Q : O \times A \rightarrow \mathbb{R}$ and a Q-function $U : X \times A \rightarrow \mathbb{R}$. The hidden state dependent Q-function U serves as a target in the temporal difference learning update. It is the tabular version of the modified DQN algorithm used in [Topin et al. \(2021\)](#). Indeed, when learning deterministic memoryless policies for IBMDPs, [Topin et al. \(2021\)](#) were using RL algorithms corresponding to asymmetric DQN or asymmetric PPO from [Baisero et al. \(2022\)](#); [Baisero & Amato \(2022\)](#) before those were published. We provide a reproducibility study in appendix 9 of [Topin et al. \(2021\)](#). In appendix 12, we describe tabular asymmetric variants such as asymmetric Q-learning, which learns an observable Q-function $Q : O \times A \rightarrow \mathbb{R}$ using a hidden-state target $U : X \times A \rightarrow \mathbb{R}$. This mirrors the modified DQN approach used in [Topin et al. \(2021\)](#), which effectively applied asymmetric DQN/PPO before these frameworks were formalized. We also provide a reproducibility study of [Topin et al. \(2021\)](#) in appendix 9. Until recently, the benefits of asymmetric RL over standard RL was only shown empirically and only for history-dependent models. The work of [Lambrechts et al. \(2025b\)](#) proves that some asymmetric RL algorithms should in theory learn better history-dependent or memoryless policies for POMDPs. This is exactly what we wish for. However, those algorithms are not practical because they require estimations of the asymptotic occupancy distribution $P^\pi((s, o')|o)$ (definition 8) for candidate policies which in turn would require to perform costly Monte-Carlo estimations. We leave it to future work to use those algorithms that combine asymmetric RL and estimation of future visitation frequencies since those results are contemporary to the writing of this manuscript.

5 Methodology

The goal of this section is to check if the aforementioned approach can consistently retrieve optimal decision tree policies for a simple grid world MDP (figure 1a). In particular, we use reinforcement learning to train decision tree policies for MDPs by seeking deterministic memoryless policies that optimize the RL objective in POIBMDPs (figure 1b and section 4).

5.1 Computing some decision tree policies

To assess the performance of reinforcement learning, for different trade-off reward ζ (definition 5), we identify the deterministic memoryless policies that maximize the RL objective (definition 2). Each of those policies correspond to one of the decision tree policies for the grid world downstream MDP illustrated in figure 2: (i) a depth-0 tree equivalent to always taking the same downstream action ($\pi_{\mathcal{T}_0}$), (ii) a depth-1 tree equivalent alternating between an IGA and a downstream action ($\pi_{\mathcal{T}_1}$), (iii) an unbalanced depth-2 tree that sometimes takes two IGAs then a downstream action and sometimes an IGA then a downstream action ($\pi_{\mathcal{T}_u}$), (iv) a depth-2 tree that alternates between taking two IGAs and a downstream action ($\pi_{\mathcal{T}_2}$), or (v) an infinite ‘tree’ that only takes IGAs. We will particularly focus on trying to retrieve the depth-1 decision tree that is the most interpretable—smallest number of nodes and shallowest—tree taking optimal actions. Because from [Singh et al. \(1994\)](#) we know that for POMDPs, stochastic memoryless policies can sometimes get better expected discounted rewards than deterministic memoryless policies, we also compute the value of the stochastic policy that randomly alternates between two downstream actions: \rightarrow and \downarrow . Taking those two downstream actions always lead to the goal state in expectation (figure 1a). Because we know all the downstream states, all the observations, all the actions, all the rewards and all the transitions of our POIBMDPs, we can compute the RL objective values of those different deterministic memoryless policies exactly given ζ and γ a discount factor. We plot, in figure 2, the RL objective values of the decision tree policies as functions of ζ when we fix $\gamma = 0.99$ (standard choice of discount in practice [Sutton & Barto \(1998\)](#)). Despite objective values being very similar for the

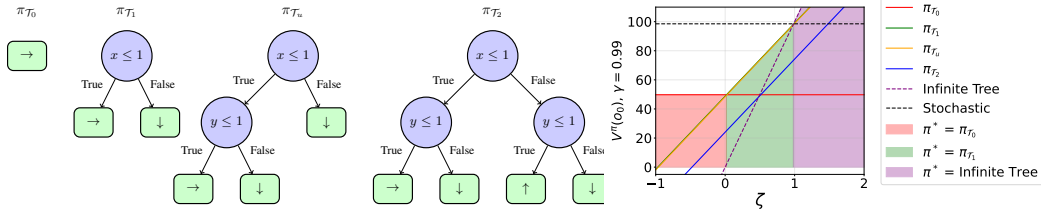


Figure 2: Decision tree policies and their RL objective values (definition 2) as functions of the rewards ζ in POIBMDPs associated to the grid world MDP (figure 1b). Shaded areas on the right plot indicate which deterministic memoryless policy is optimal depending on ζ . Recall that ζ can be seen as a reward that encourages the deterministic memoryless policy to take information-gathering actions, i.e. the reward that trades off interpretability of the corresponding decision tree and its cumulative reward.

depth-1 and unbalanced depth-2 tree, we now know from the green shaded area that a depth-1 tree is optimal when $0 < \zeta < 1$ in the POIBMDP. Interestingly, two challenges of learning in POMDPs described in Singh et al. (1994) are visible in figure 2. First, there is a whole range of ζ values for which the optimal memoryless policy is stochastic. Second, for e.g. $\zeta = 0.5$, while the optimal deterministic memoryless policy corresponds to a depth-1 tree, the value of the (PO)IBMDP state $(s_2, o_0) = (1.5, 1.5, 0, 2, 0, 2)$, i.e. the agent current position is upper right cell in the downstream MDP but it has not information about it, is not maximized by the optimal memoryless policy that will take 1 information-gathering action in between each downstream action (figure 2), but by the sub-optimal policy that always goes down. Next, we present the specific experimental setup that we use in the remaining of the paper.

5.2 Experimental setup

All our code is open source¹. The downstream MDP for which we want to learn decision tree policies with reinforcement learning is the grid world MDP (figure 1a). We then get 100 associated POIBMDPs following figure 1b with ζ chosen uniformly among 100 different in $]0, 1[$.

We will evaluate two types of reinforcement learning algorithm. First we use standard tabular RL algorithms, namely Q-learning, Sarsa, and vanilla policy gradient on a softmax policy Watkins & Dayan (1992); Sutton & Barto (1998); Jaakkola et al. (1994), to learn deterministic memoryless policies in POIBMDPs by simply replacing the current state in the algorithm descriptions by the current observation. 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 sought depth-1 decision tree policy (figure 2), but for completeness we will see what trees are obtained after greddification of the stochastic policies.

Second, we also use the more specialised asymmetric Q-learning, asymmetric Sarsa, and asymmetric policy gradient (algorithm 4, section 4.2). Each algorithm is trained until convergence on each POIBMDP, and each one of those runs is repeated 100 times. For all baselines we use, when applicable, exploration rates $\epsilon = 0.3$ and learning rates $\alpha = 0.1$. All the training curves are presented in appendix 11.

We will consider two metrics. We consider the distribution of the learned trees over the 100 training seeds. Indeed, since for every POIBMDP we run each algorithm 100 times, at the end of training we get 100 deterministic memoryless policies, from which we can extract the equivalent 100 decision tree policies using algorithm 1 and we can count which one have e.g. a depth of 1. This helps understand which trees RL algorithms tend to learn as a function of the trade-off reward ζ .

¹<https://anonymous.4open.science/r/poibmdps-5BFE/>

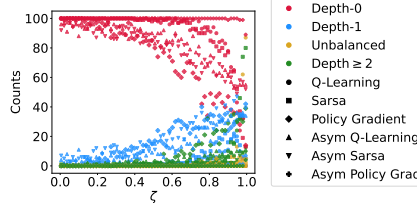


Figure 3: Distributions of final tree policies learned across the 100 seeds. For each ζ value, there are four colored points. Each point represent the share of depth-0 trees (red), depth-1 trees (green), unbalanced depth-2 trees (orange) and depth-2 trees (blue).

5.3 Can reinforcement learning find the optimal deterministic memoryless POIBMDP policies?

In figure 13, we plot the distributions of the final learned trees over the 100 random seeds in function of ζ from the above runs. For example, in figure 13, in the top left plot, when learning 100 times in a POIBMDP with $\zeta = 0.5$, Q-learning returned almost 100 times a depth-0 tree. Again, on none of those subplots do we see a high rate of learned depth-1 trees for $\zeta \in]0, 1[$. It is alerting that the most frequent learned trees are the depth-0 trees for $\zeta \in]0, 1[$ because such trees are way more sub-optimal than e.g. the depth-2 unbalanced trees (figure 2). 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 downstream action. Furthermore, in appendix 11 we show that for POIBMDPs with $\zeta \in [-1, 0]$ and $\zeta \in [1, 2]$, baselines consistently learn the optimal policies—a depth-0 tree and an infinite tree respectively—which is concerning as it means that RL seem to find only trivial policies. On the positive side, we observe that asymmetric versions of Q-learning and Sarsa have found the optimal deterministic memoryless policy—the depth-1 decision tree—more frequently throughout the optimality range $]0, 1[$, than their symmetric counter-parts for $\zeta \in]0, 1[$. Next, we quantify how difficult it is to do RL to learn memoryless policies in POIBMDPs as opposed to standard Markovian policies.

5.4 How difficult is it to learn in POIBMDPs?

In this section we run the same (asymmetric) reinforcement learning algorithms to learn standard Markovian policies in MDPs (definition 1) or IBMDPs (definition 5), or deterministic memoryless policies in POIBMDPs (definition 7).

In order to see how difficult each of these three problems is, we can run a *great* number of experiments for each problem and compare solving rates. To make solving rates comparable we consider a unique instance for each of those problems. Problem 1 is learning one of the optimal standard Markovian deterministic policy ($\pi : S \rightarrow A$) from figure 1a for the grid world MDP with $\gamma = 0.99$. Problem 2 is learning one of the optimal standard Markovian deterministic policy ($\pi : S \times O \rightarrow A \cup A_{info}$) for the IBMDP from figure 1b with $\gamma = 0.99$ and $\zeta = 0.5$. Problem 3 is what has been done in the previous section to learn deterministic memoryless policies ($\pi : O \rightarrow A \cup A_{info}$) where in addition of fixing $\gamma = 0.99$ we also fix $\zeta = 0.5$.

We use the six (asymmetric) RL algorithms from the previous section and try a wide set of hyperparameters and additional learning tricks (optimistic Q-function, eligibility traces, entropy regularization and ϵ -decay, all are described in Sutton & Barto (1998)). The complete detailed lists of hyperparameters are given in the appendix 12 and a summary is given in table 6. Furthermore, the careful reader might notice that there is no point running asymmetric RL on MDPs or IBMDPs when the problem does not require partial observability. Hence, we only run asymmetric RL for POIBMDPs and otherwise run all other RL algorithms and all problems.

Each unique hyperparameter combination for a given algorithm on a given problem is run 10 times on 1 million learning steps to get standard errors. For example, for asymmetric Sarsa, we run a total

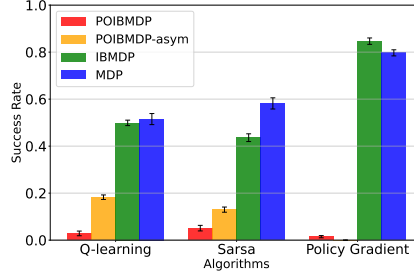


Figure 4: Success rates of different (asymmetric) RL algorithms over thousands of runs when applied to learning deterministic memoryless policies in a POIBMDP or learning deterministic policies in associated MDP and IBMDP.

of $10 \times 768 = 7680$ experiments for learning deterministic memoryless policies for a POIBMDP. To get a success rate, we can simply divide the number of learned optimal depth-1 tree by 768 (recall that for $\gamma = 0.99$ and $\zeta = 0.5$, the optimal policy is a depth-1 tree (e.g. figure 2)).

The key observations from figure 4 is that reinforcement learning a deterministic memoryless policy in a POIBMDP, is way harder than learning a standard Markovian policy. For example, Q-learning finds the optimal solution in only 3% of the experiments while the same algorithms to optimize the standard RL objective (definition 2) in an MDP or IBMDP found the optimal solutions 50% of the time. 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 originally developed in [Topin et al. \(2021\)](#) seems way too 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 here, 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: how to instantiate the (PO)IBMDP action space such that internal nodes in resulting trees are useful for predictions? Next, we further support that partial observability is the main limitation for reinforcement learning to train decision tree policies for MDPs.

5.5 Are they downstream MDPs for which partial observability is not a limitation?

In this section, we show that for a special class of POIBMDPs, reinforcement learning can learn optimal deterministic memoryless policies w.r.t the RL objective, i.e. we can do direct decision tree policy learning for MDPs. This class of POIBMDPs are those for which downstream MDPs have uniform transitions, i.e. $T(s, a, s') = \frac{1}{|S|}$ (definitions 1 and 5). Such downstream MDPs include classification tasks formulated as MDPs. This implies that learning deterministic memoryless policies in POIBMDPs where the downstream MDP encodes a classification task is equivalent to doing supervised learning of a decision tree (figure 5). This is exactly what is done in e.g. [Kohler et al. \(2025a\)](#). In figure 5 we give an example of such downstream MDPs for a classification task with 4 data in the training set and 2 classes: $\mathcal{X} = \{(0.5, 0.5), (0.5, 1.5), (1.5, 1.5), (1.5, 0.5)\}$ and $y = \{0, 0, 1, 1\}$.

In appendix 13, we show that POIBMDPs associated to downstream MDPs with uniform transitions are themselves standard MDP (definition 1). This means that in principle, standard RL algorithms like Q-learning, should work as well as for any MDP. If RL does work for such fully observable POIBMDPs, this would mean that the difficulty of direct learning of decision tree policies for *any* MDP using POIBMDPs, exhibited in sections, is most likely due to the partial observability. This is exactly what we check next. We use the same direct approach to learn decision tree policies as in previous sections, except that now the downstream MDP is a classification task and not a sequential decision making task like reaching a goal in a grid world.

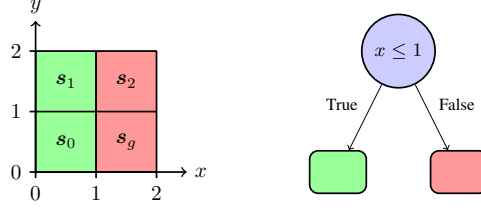


Figure 5: In this MDP, there are four data to which to assign either a green or red label. On the right, there is the unique optimal depth-1 tree for this particular MDP. This depth-1 tree also maximizes the accuracy on the corresponding classification task.

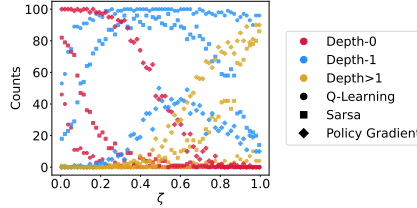


Figure 6: We reproduce the same plot as in figure 13 for POIBMDPs associated to a downstream MDP encoding a classification task. Each colored dot is the number of final learned trees with a specific structure for a given ζ .

417 We construct POIBMDPs for the classification task from figure 5, with $\gamma = 0.99$, 200 values of
 418 $\zeta \in [0, 1]$ and IGAs $x \leq 1$ and $y \leq 1$. Since those POIBMDPs are MDPs, we do not need to
 419 analyze asymmetric RL baselines. We see on figure 6 that compared to general POIBMDPs from
 420 previous sections, RL can be used to consistently learn optimal deterministic memoryless policies
 421 $O := A \cup A_{info}$. Such policies are equivalent to decision tree classifiers.

422 6 Discussion

423 In this paper, we were interested in algorithms that can learn decision tree policies that directly opti-
 424 mize some trade-off of interpretability and performance in MDPs without relying on an oracle or
 425 expert policy. Starting from the IBMDP formulation from Topin et al. (2021), we have shown that
 426 direct learning of decision tree policies for MDPs can be framed as learning deterministic mem-
 427 oryless policies in POMDPs that we called POIBMDPs. By bridging the gap with the POMDP
 428 literature, we conjectured that partial observability is the main limitation for reinforcement learning
 429 of decision tree policies for MDPs. We then supported this conjecture by benchmarking different
 430 reinforcement learning algorithms on carefully crafted problems for which we knew the exact opti-
 431 mal decision tree policies. Across our experiments, we found that only when partial observability is
 432 absent from the learning task can good decision tree policies be trained. Attempting to overcome the
 433 partial observability challenges highlighted so far seems like a bad research direction. Indeed, while
 434 algorithms tailored specifically for the problem of learning deterministic memoryless policies for
 435 POIBMDPs might exist, imitation learning works well in practice and has been the state-of-the-art
 436 for interpretable sequential decision making for a while. Furthermore there are other limitations that
 437 we did not cover in the framework of Topin et al. (2021) such as how to choose good candidates
 438 information-gathering actions or simply how to choose ζ for a target interpretability-performance
 439 trade-off. Finally, while we focused on non-parametric tree learning assuming RL algorithms should
 440 naturally trade off interpretability and performance through the reward signal ζ for adding nodes to
 441 the decision tree policy, another future research avenue could be to develop better algorithms train-
 442 ing parametric trees. Indeed parametric tree policies, on the other hand, can be computed with
 443 reinforcement learning directly in the downstream MDP. However such RL algorithms for paramet-
 444 ric decision tree policies Silva et al. (2020); Vos & Verwer (2024); Marton et al. (2025) require to

re-train a policy entirely for each desired level of interpretability, i.e. each unique tree structure. Future research in this direction should focus on algorithms for parametric tree policies that can re-use samples from one tree learning to train a different tree structure more efficiently. This would reduce the required quantity of a priori knowledge on the decision tree policy structure.

Broader Impact Statement

Our work put great emphasis on covering existing work, open sourcing code, and being transparent about limitations. We hope that the impact of our work is going to be positive for society through advancing research in interpretable machine learning.

Acknowledgments

Use unnumbered third level headings for the acknowledgments. All acknowledgments, including those to funding agencies, go at the end of the paper. Only add this information once your submission is accepted and deanonymized. The acknowledgments do not count towards the 8–12 page limit.

References

- Mauricio Araya-López, Vincent Thomas, Olivier Buffet, and François Charpillet. A Closer Look at MOMDPs. In *Proceedings of the 22nd International Conference on Tools with Artificial Intelligence*, Proceedings of the 22nd International Conference on Tools with Artificial Intelligence, Arras, France, October 2010. IEEE. URL <https://inria.hal.science/inria-00535559>.
- Akanksha Atrey, Kaleigh Clary, and David Jensen. Exploratory not explanatory: Counterfactual analysis of saliency maps for deep reinforcement learning. In *International Conference on Learning Representations*, 2020. URL <https://openreview.net/forum?id=rkl3mlBFDB>.
- Andrea Baisero and Christopher Amato. Unbiased asymmetric reinforcement learning under partial observability. In *Proceedings of the 21st International Conference on Autonomous Agents and Multiagent Systems*, AAMAS '22, pp. 44–52, Richland, SC, 2022. International Foundation for Autonomous Agents and Multiagent Systems. ISBN 9781450392136.
- Andrea Baisero, Brett Daley, and Christopher Amato. Asymmetric DQN for partially observable reinforcement learning. In James Cussens and Kun Zhang (eds.), *Proceedings of the Thirty-Eighth Conference on Uncertainty in Artificial Intelligence*, volume 180 of *Proceedings of Machine Learning Research*, pp. 107–117. PMLR, 01–05 Aug 2022. URL <https://proceedings.mlr.press/v180/baisero22a.html>.
- Andrew G. Barto, Richard S. Sutton, and Charles W. Anderson. Neuronlike adaptive elements that can solve difficult learning control problems. *IEEE Transactions on Systems, Man, and Cybernetics*, SMC-13(5):834–846, 1983. DOI: 10.1109/TSMC.1983.6313077.
- Osbert Bastani, Yewen Pu, and Armando Solar-Lezama. Verifiable reinforcement learning via policy extraction. 2018.
- Richard Bellman. *Dynamic Programming*. 1957.
- Dimitris Bertsimas and Jack Dunn. Optimal classification trees. *Machine Learning*, 106:1039–1082, 2017.
- Guy Blanc, Jane Lange, and Li-Yang Tan. Learning stochastic decision trees. *CoRR*, abs/2105.03594, 2021. URL <https://arxiv.org/abs/2105.03594>.
- L Breiman, JH Friedman, R Olshen, and CJ Stone. *Classification and Regression Trees*. Wadsworth, 1984.

- 487 Samuel Ping-Man Choi, Nevin Lianwen Zhang, and Dit-Yan Yeung. Solving hidden-mode markov
488 decision problems. In Thomas S. Richardson and Tommi S. Jaakkola (eds.), *Proceedings of*
489 *the Eighth International Workshop on Artificial Intelligence and Statistics*, volume R3 of *Pro-*
490 *ceedings of Machine Learning Research*, pp. 49–56. PMLR, 04–07 Jan 2001. URL [https://](https://proceedings.mlr.press/r3/choi01a.html)
491 proceedings.mlr.press/r3/choi01a.html. Reissued by PMLR on 31 March
492 2021.
- 493 Emir Demirovic, Anna Lukina, Emmanuel Hebrard, Jeffrey Chan, James Bailey, Christopher
494 Leckie, Kotagiri Ramamohanarao, and Peter J. Stuckey. Murtree: Optimal decision trees via
495 dynamic programming and search. *Journal of Machine Learning Research*, 23(26):1–47, 2022.
496 URL <http://jmlr.org/papers/v23/20-520.html>.
- 497 Emir Demirović, Emmanuel Hebrard, and Louis Jean. Blossom: an anytime algorithm for com-
498 puting optimal decision trees. *Proceedings of the 40th International Conference on Machine*
499 *Learning*, 202:7533–7562, 23–29 Jul 2023. URL [https://proceedings.mlr.press/](https://proceedings.mlr.press/v202/demirovic23a.html)
500 [v202/demirovic23a.html](https://proceedings.mlr.press/v202/demirovic23a.html).
- 501 Claire Glanois, Paul Weng, Matthieu Zimmer, Dong Li, Tianpei Yang, Jianye Hao, and Wulong Liu.
502 A survey on interpretable reinforcement learning. *Machine Learning*, pp. 1–44, 2024.
- 503 Tommi Jaakkola, Satinder P. Singh, and Michael I. Jordan. Reinforcement learning algorithm for
504 partially observable markov decision problems. In *Proceedings of the 8th International Confer-*
505 *ence on Neural Information Processing Systems*, NIPS’94, pp. 345–352, Cambridge, MA, USA,
506 1994. MIT Press.
- 507 Hector Kohler, Riad Akrou, and Philippe Preux. Breiman meets bellman: Non-greedy decision
508 trees with mdps. In *Proceedings of the 31st ACM SIGKDD Conference on Knowledge Discovery*
509 *and Data Mining V.2*, KDD ’25, pp. 1207–1218, New York, NY, USA, 2025a. Association for
510 Computing Machinery. ISBN 9798400714542. DOI: 10.1145/3711896.3736868. URL [https://](https://doi.org/10.1145/3711896.3736868)
511 doi.org/10.1145/3711896.3736868.
- 512 Hector Kohler, Waris Radji, Quentin Delfosse, Riad Akrou, and Philippe Preux. Evaluating inter-
513 pretable reinforcement learning by distilling policies into programs. In *RLC 2025 Workshop on*
514 *Programmatic Reinforcement Learning*, 2025b. URL [https://openreview.net/forum?](https://openreview.net/forum?id=n1CfzixauT)
515 [id=n1CfzixauT](https://openreview.net/forum?id=n1CfzixauT).
- 516 Gaspard Lambrechts, Adrien Bolland, and Damien Ernst. Informed POMDP: Leveraging additional
517 information in model-based RL. *Reinforcement Learning Journal*, 2:763–784, 2025a.
- 518 Gaspard Lambrechts, Damien Ernst, and Aditya Mahajan. A theoretical justification for asymmetric
519 actor-critic algorithms. In *Forty-second International Conference on Machine Learning*, 2025b.
520 URL <https://openreview.net/forum?id=FlyANMCnAn>.
- 521 Zachary C. Lipton. The mythos of model interpretability: In machine learning, the concept of
522 interpretability is both important and slippery. *Queue*, 16(3):31–57, 2018.
- 523 Michael L. Littman. Memoryless policies: theoretical limitations and practical results. In *Proceed-*
524 *ings of the Third International Conference on Simulation of Adaptive Behavior: From Animals to*
525 *Animats 3: From Animals to Animats 3*, SAB94, pp. 238–245, Cambridge, MA, USA, 1994. MIT
526 Press. ISBN 0262531224.
- 527 John Loch and Satinder P. Singh. Using eligibility traces to find the best memoryless policy in
528 partially observable markov decision processes. In *Proceedings of the Fifteenth International*
529 *Conference on Machine Learning*, ICML ’98, pp. 323–331, San Francisco, CA, USA, 1998.
530 Morgan Kaufmann Publishers Inc. ISBN 1558605568.
- 531 Scott M. Lundberg and Su-In Lee. A unified approach to interpreting model predictions. In *Proceed-*
532 *ings of the 31st International Conference on Neural Information Processing Systems*, NIPS’17,
533 pp. 4768–4777, Red Hook, NY, USA, 2017. Curran Associates Inc. ISBN 9781510860964.

- 534 Yishay Mansour, Michal Moshkovitz, and Cynthia Rudin. There is no accuracy-interpretability
535 tradeoff in reinforcement learning for mazes, 2022. URL [https://arxiv.org/abs/2206.](https://arxiv.org/abs/2206.04266)
536 [04266](https://arxiv.org/abs/2206.04266).
- 537 Sascha Marton, Tim Grams, Florian Vogt, Stefan Lüdtkke, Christian Bartelt, and Heiner Stucken-
538 schmidt. Mitigating information loss in tree-based reinforcement learning via direct optimization.
539 2025. URL <https://openreview.net/forum?id=qpXctF2aLZ>.
- 540 Stephanie Milani, Nicholay Topin, Manuela Veloso, and Fei Fang. Explainable reinforcement learn-
541 ing: A survey and comparative review. *ACM Comput. Surv.*, 56(7), April 2024. ISSN 0360-0300.
542 DOI: 10.1145/3616864. URL <https://doi.org/10.1145/3616864>.
- 543 Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Belle-
544 mare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level
545 control through deep reinforcement learning. *nature*, 518(7540):529–533, 2015.
- 546 Sreerama Murthy and Steven Salzberg. Lookahead and pathology in decision tree induction. In
547 *Proceedings of the 14th International Joint Conference on Artificial Intelligence - Volume 2*,
548 IJCAI’95, pp. 1025–1031, San Francisco, CA, USA, 1995. Morgan Kaufmann Publishers Inc.
549 ISBN 1558603638.
- 550 F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Pretten-
551 hofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and
552 E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*,
553 12:2825–2830, 2011.
- 554 Lerrel Pinto, Marcin Andrychowicz, Peter Welinder, Wojciech Zaremba, and Pieter Abbeel. Asym-
555 metric actor critic for image-based robot learning, 2017. URL [https://arxiv.org/abs/](https://arxiv.org/abs/1710.06542)
556 [1710.06542](https://arxiv.org/abs/1710.06542).
- 557 Nikaash Puri, Sukriti Verma, Piyush Gupta, Dhruv Kayastha, Shripad Deshmukh, Balaji Krishna-
558 murthy, and Sameer Singh. Explain your move: Understanding agent actions using specific and
559 relevant feature attribution. In *International Conference on Learning Representations*, 2020. URL
560 <https://openreview.net/forum?id=SJgzLkBKPB>.
- 561 Antonin Raffin, Ashley Hill, Adam Gleave, Anssi Kanervisto, Maximilian Ernestus, and Noah Dor-
562 mann. Stable-baselines3: Reliable reinforcement learning implementations. *Journal of Machine*
563 *Learning Research*, 22(268):1–8, 2021.
- 564 Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. "why should i trust you?": Explaining
565 the predictions of any classifier. pp. 1135–1144, 2016. DOI: 10.1145/2939672.2939778. URL
566 <https://doi.org/10.1145/2939672.2939778>.
- 567 Stéphane Ross, Geoffrey J. Gordon, and J. Andrew Bagnell. A reduction of imitation learning and
568 structured prediction to no-regret online learning. 2010.
- 569 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy
570 optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- 571 Wenjie Shi, Gao Huang, Shiji Song, Zhuoyuan Wang, Tingyu Lin, and Cheng Wu. Self-supervised
572 discovering of interpretable features for reinforcement learning. *IEEE Transactions on Pat-*
573 *tern Analysis and Machine Intelligence*, 44(5):2712–2724, 2022. DOI: 10.1109/TPAMI.2020.
574 3037898.
- 575 Olivier Sigaud and Olivier Buffet. *Partially Observable Markov Decision Processes*, chapter 7,
576 pp. 185–228. John Wiley Sons, Ltd, 2013. ISBN 9781118557426. DOI: [https://doi.org/](https://doi.org/10.1002/9781118557426.ch7)
577 [10.1002/9781118557426.ch7](https://doi.org/10.1002/9781118557426.ch7). URL [https://onlinelibrary.wiley.com/doi/abs/](https://onlinelibrary.wiley.com/doi/abs/10.1002/9781118557426.ch7)
578 [10.1002/9781118557426.ch7](https://onlinelibrary.wiley.com/doi/abs/10.1002/9781118557426.ch7).

- Andrew Silva, Matthew Gombolay, Taylor Killian, Ivan Jimenez, and Sung-Hyun Son. Optimization methods for interpretable differentiable decision trees applied to reinforcement learning. In Silvia Chiappa and Roberto Calandra (eds.), *Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics*, volume 108 of *Proceedings of Machine Learning Research*, pp. 1855–1865. PMLR, 26–28 Aug 2020. URL <https://proceedings.mlr.press/v108/silva20a.html>.
- Satinder P. Singh, Tommi S. Jaakkola, and Michael I. Jordan. Learning without state-estimation in partially observable markovian decision processes. In *Proceedings of the Eleventh International Conference on International Conference on Machine Learning*, ICML’94, pp. 284–292, San Francisco, CA, USA, 1994. Morgan Kaufmann Publishers Inc. ISBN 1558603352.
- Edward J. Sondik. The optimal control of partially observable markov processes over the infinite horizon: Discounted costs. *Operations Research*, 26(2):282–304, 1978. ISSN 0030364X, 15265463. URL <http://www.jstor.org/stable/169635>.
- Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning: An Introduction*. The MIT Press, Cambridge, MA, 1998.
- Richard S Sutton, David McAllester, Satinder Singh, and Yishay Mansour. Policy gradient methods for reinforcement learning with function approximation. In S. Solla, T. Leen, and K. Müller (eds.), *Advances in Neural Information Processing Systems*, volume 12. MIT Press, 1999. URL https://proceedings.neurips.cc/paper_files/paper/1999/file/464d828b85b0bed98e80ade0a5c43b0f-Paper.pdf.
- Nicholay Topin, Stephanie Milani, Fei Fang, and Manuela Veloso. Iterative bounding mdps: Learning interpretable policies via non-interpretable methods. *Proceedings of the AAAI Conference on Artificial Intelligence*, 35:9923–9931, 2021.
- Mark Towers, Ariel Kwiatkowski, Jordan Terry, John U Balis, Gianluca De Cola, Tristan Deleu, Manuel Goulão, Andreas Kallinteris, Markus Krimmel, Arjun KG, et al. Gymnasium: A standard interface for reinforcement learning environments. *arXiv preprint arXiv:2407.17032*, 2024.
- Jacobus van der Linden, Mathijs de Weerdt, and Emir Demirović. Necessary and sufficient conditions for optimal decision trees using dynamic programming. *Advances in Neural Information Processing Systems*, 36:9173–9212, 2023.
- Abhinav Verma, Vijayaraghavan Murali, Rishabh Singh, Pushmeet Kohli, and Swarat Chaudhuri. Programmatically interpretable reinforcement learning. pp. 5045–5054, 2018.
- Sicco Verwer and Yingqian Zhang. Learning optimal classification trees using a binary linear program formulation. *Proceedings of the AAAI conference on artificial intelligence*, 33:1625–1632, 2019.
- Daniël Vos and 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. ISBN 978-1-956792-03-4. DOI: 10.24963/ijcai.2023/606. URL <https://doi.org/10.24963/ijcai.2023/606>.
- Daniël Vos and Sicco Verwer. Optimizing interpretable decision tree policies for reinforcement learning. 2024. URL <https://arxiv.org/abs/2408.11632>.
- Christopher JCH Watkins and Peter Dayan. Q-learning. *Machine learning*, 8(3):279–292, 1992.
- Haoze Wu, Omri Isac, Aleksandar Zeljić, Teruhiro Tagomori, Matthew Daggitt, Wen Kokke, Idan Refaeli, Guy Amir, Kyle Julian, Shahaf Bassan, Pei Huang, Ori Lahav, Min Wu, Min Zhang, Ekaterina Komendantskaya, Guy Katz, and Clark Barrett. Marabou 2.0: A versatile formal analyzer of neural networks, 2024. URL <https://arxiv.org/abs/2401.14461>.

Data: Deterministic partially observable policy π_{po} for IBMDP

$\mathcal{M}_{IB} \equiv \langle S \times O, A \cup A_{info}, (R, \zeta), (T_{info}, T, T_0) \rangle$ and IBMDP observation

$\mathbf{o} = (L'_1, U'_1, \dots, L'_p, U'_p)$

Result: Decision tree policy $\pi_{\mathcal{T}}$ for MDP $\mathcal{M} \equiv \langle S, A, R, T, T_0 \rangle$

Function Subtree_From_Policy(\mathbf{o}, π_{po}):

```

   $a \leftarrow \pi_{po}(\mathbf{o})$ 
  if  $a$  is a downstream action then
    | return Leaf_Node(action:  $a$ )
  end
  else
     $\langle i, v \rangle \leftarrow a$ 
     $\mathbf{o}_L \leftarrow \mathbf{o}; \quad \mathbf{o}_R \leftarrow \mathbf{o}$ 
     $\mathbf{o}_L \leftarrow (L'_1, U'_1, \dots, L'_j, v, \dots, L'_p, U'_p); \quad \mathbf{o}_R \leftarrow (L'_1, U'_1, \dots, v, U'_j, \dots, L'_p, U'_p)$ 
     $child_L \leftarrow \text{Subtree\_From\_Policy}(\mathbf{o}_L, \pi_{po})$ 
     $child_R \leftarrow \text{Subtree\_From\_Policy}(\mathbf{o}_R, \pi_{po})$ 
    return Internal_Node(feature:  $i$ , value:  $v$ , children: ( $child_L, child_R$ ))
  end

```

Algorithm 1: Extract a decision tree policy (algorithm 1 from [Topin et al. \(2021\)](#))

Supplementary Materials

The following content was not necessarily subject to peer review.

7 From a policy to a tree

8 Imitation learning: a baseline for indirect decision tree policy learning

In this section we present decision tree policies of this manuscript obtained using Dagger or VIPER [Bastani et al. \(2018\)](#); [Verma et al. \(2018\)](#) after learning an expert Q-function for the grid world MDP. Recall the optimal policies for the grid world, taking the green actions in each state in figure 1a. Among the optimal policies, the ones that go left or up in the goal state can be problematic for imitation learning algorithms. Indeed, we know that for this grid world MDP there exists decision tree policies with a very good interpretability-performance trade-off: depth-1 decision trees that are optimal w.r.t. the RL objective. One could even say that those trees have the *optimal* interpretability-performance trade-off because they are the shortest trees that are optimal w.r.t. the RL objective.

In figure 7, we present a depth-1 decision tree policy that is optimal w.r.t. the RL objective and a depth-1 tree that is sub-optimal. The other optimal depth-1 tree is to go right when $y \leq 1$ and down otherwise.

Now a fair question is: can Dagger or VIPER learn such an optimal depth-1 tree given access to an expert optimal policy from figure 1a?

We start by running the standard Q-learning algorithm with $\epsilon = 0.3$, $\alpha = 0.1$ over 10,000 time steps. While for Q-learning, Sutton and Barto break ties by index value in their book [Sutton & Barto \(1998\)](#) (the greedy action is the argmax action with smallest index), we show that the choice of tie-breaking greatly influences the performance of subsequent imitation learning algorithms. Indeed, depending on how actions are ordered in practice, Q-learning may be biased toward some optimal policies rather than others. While this does not matter for one who just wants to find an optimal policy, in our example of finding the optimal depth-1 decision tree policy, it matters *a lot*.

In the left plot of figure 8, we see that Q-learning, independently of how ties are broken, consistently converges to an optimal policy over 100 runs (random seeds). However, in the right plot of figure 8, where we plot the proportion over 100 runs of optimal decision trees returned by Dagger or VIPER

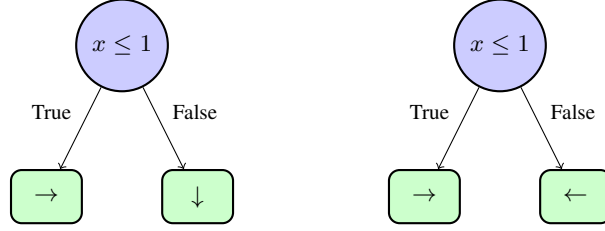


Figure 7: Left, an optimal depth-1 decision tree policy for the grid world MDP from figure 1a. On the right, a sub-optimal depth-1 decision tree policy.

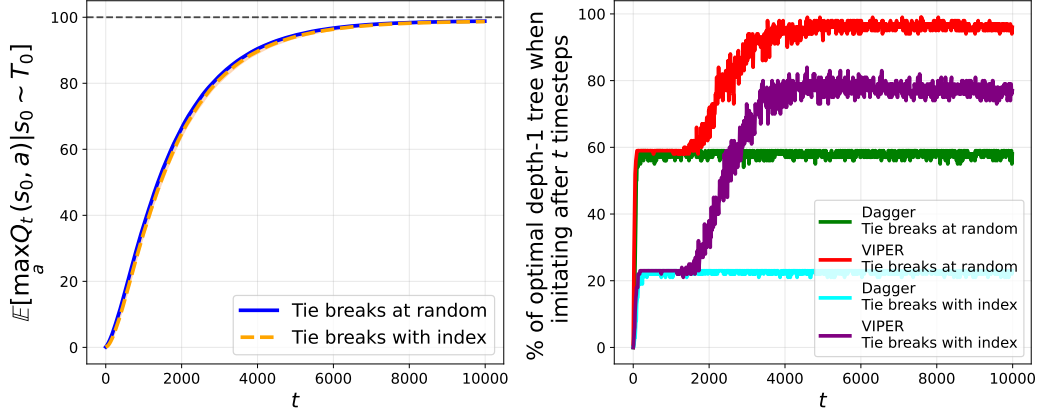


Figure 8: Left, sample complexity curve of Q-learning with default hyperparameters on the 2×2 grid world MDP over 100 random seeds. Right, performance of indirect interpretable methods when imitating the greedy policy with a tree at different Q-learning stages.

at different stages of Q-learning, we observe that imitating the optimal policy obtained by breaking ties at random consistently yields more optimal trees than breaking ties by indices. What actually happens is that the most likely output of Q-learning when ties are broken by indices is the optimal policy that goes left in the goal state, which cannot be perfectly represented by a depth-1 decision tree, because there are three different actions taken and a binary tree of depth $D = 1$ can only map to $2^D = 2$ labels.

This short experiment shows that imitation learning approaches can sometimes be very bad at learning decision tree policies with good interpretability-performance trade-offs for very simple MDPs. Despite VIPER almost always finding the optimal depth-1 decision tree policy in terms of the RL objective when ties are broken at random, we have shed light on the sub-optimality of indirect approaches such as imitation learning. This motivates the study of direct approaches to directly search for policies with good interpretability-performance trade-offs with respect to the original RL objective.

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

We attempt to reproduce the results from [Topin et al. \(2021\)](#) in which authors compare direct and indirect learning of decision tree policies of depth at most 2 for the CartPole MDP [Barto et al. \(1983\)](#). In the original paper, the authors find that both direct and indirect learning yields decision tree policies with similar RL objective values (definition 2) for the CartPole. On the other hand, we find that, imitation learning, despite not directly optimizing the RL objective for CartPole, outperforms deep RL which directly optimizes a trade-off of the standard RL objective and interpretability.

Authors of [Topin et al. \(2021\)](#) use two deep reinforcement learning baselines to which they apply some modifications in order to learn memoryless policies. Authors modify the standard DQN ? to return a memoryless 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. We present the modifications in algorithm 2. Similar modifications are applied to the standard PPO [Schulman et al. \(2017\)](#) that we present in the appendix (algorithm 3). In the modified PPO, neural network policy $O \rightarrow A \cup A_{info}$ is trained using a neural network value function $S \times O \rightarrow A \cup A_{info}$ as a critic.

Those two variants of DQN and PPO have first been introduced in [Pinto et al. \(2017\)](#) for robotic tasks with partially observable components, under the name “asymmetric” actor-critic. Asymmetric RL algorithms that have policy and value estimates using different information from a POMDP [Sondik \(1978\)](#); [Sigaud & Buffet \(2013\)](#) were later studied theoretically to solve POMDPs in Baisero’s work [Baisero et al. \(2022\)](#); [Baisero & Amato \(2022\)](#). The connections from Deep RL in IBMDPs for objective is absent from [Topin et al. \(2021\)](#) and we defer their connections to direct interpretable reinforcement learning to the next chapter as our primary goal is to reproduce [Topin et al. \(2021\)](#) as is. Next, we present the precise experimental setup we use to reproduce [Topin et al. \(2021\)](#) in order to study direct deep reinforcement learning of decision tree policies for the CartPole MDP.

9.1 IBMDP formulation

Given a base MDP $\mathcal{M} \equiv \langle S, A, R, T, T_0 \rangle$ (cf. definition 1), in order to define an IBMDP $\mathcal{M}_{IB} \langle S \times O, A \cup A_{info}, (R, \zeta), (T, T_0, T_{info}) \rangle$ (cf. definition 5), the user needs to provide the set of information gathering actions A_{info} and the reward ζ for taking those. Authors of [Topin et al. \(2021\)](#) propose to parametrize the set of IGAs with $j \times q$ actions $\langle j, v_k \rangle$ with v_k depending on the current observation $\mathbf{o}_t = (L'_1, U'_1, \dots, L'_j, U'_j, \dots, L'_p, U'_p)$: $v_k = \frac{k(U'_j - L'_j)}{q+1}$. This parametric IGAs space keeps the discrete IBMDP action space at a reasonable size while providing a learning algorithm with varied IGAs to try.

For example, if we define an IBMDP with $q = 3$ for the grid world from figure 1a, the grid world action space is augmented with six IGAs. At $t = 0$, recall that $\mathbf{o}_0 = (0, 2, 0, 2)$, so if an IGA is taken, e.g. $\langle 2, v_2 \rangle$, the effective IGA is $\langle j, v_2 = \frac{k(2-0)}{3+1} \rangle = \langle 1, 2 \rangle$ which in turn effectively corresponds to an internal decision tree node $y \leq 1$. If the current state y -feature value is 0.5, then the next observation at $t = 1$ is $\mathbf{o}_1 = (0, 2, 0, 1)$. At $t = 2$ if $a_t = \langle 2, v_2 \rangle$ again, it would be effectively $\langle j, v_2 = \frac{k(1-0)}{3+1} \rangle = \langle 2, 0.5 \rangle$. This would give the next observation at $t = 2$ $\mathbf{o}_2 = (0, 2, 0, 0.5)$ and so on.

Furthermore, author propose to regularize the learned decision tree policy with a maximum depth parameter D . Unfortunately, the authors did not describe how 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 in the IBMDP, we can either give negative reward for taking D IGAs in a row, or terminate the trajectory. In practice, we could also have a state-dependent action space such that taking an IGA is not allowed after taking D IGAs in a row. The latter approach—sometimes called action masking—is not compatible with the definition of an MDP (cf. definition 1) in which all actions are available in all states. To apply the penalization approaches, one can extend the MDP states to keep track of the current tree depth. Similarly, the termination approach requires a transition function that depends on the current tree depth.

We actually find that when $q + 1$, the parameter that defines threshold values in decision tree policy nodes (cf. definition 5), is a prime number, then as a direct consequence of the *Chinese Remainder Theorem*², the current tree depth is directly encoded in the current observation \mathbf{o}_t . Hence, when $q + 1$ is prime, we can control the depth through either transitions or rewards without tracking the

²https://en.wikipedia.org/wiki/Chinese_remainder_theorem

Table 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 γ	1
Information gathering actions parameter q	2, 3
Information gathering actions rewards ζ	-0.01, 0.01
Depth control	Done signal, negative reward, none

Table 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	10^6
Random transitions before learning	10^5
Epsilon start	0.9, 0.5
Epsilon end	0.05
Exploration fraction	0.1
Optimizer	RMSprop ($\alpha = 0.95$)
Learning rate	2.5×10^{-4}
Networks architectures	[128, 128]
Networks activations	tanh(), relu()

tree depth. We use the exact same downstream MDP and associated IBMDPs for our experiments as [Topin et al. \(2021\)](#) except when mentioned otherwise.

downstream MDP The task at hand is to optimize the RL objective (definition 2) with a decision tree policy for the CartPole MDP [Barto et al. \(1983\)](#). At each time step a learning algorithm observes the cart’s position and velocity and the pole’s angle and angular velocity, and can take action to push the CartPole left or right. While the CartPole is roughly balanced, i.e., while the cart’s angle remains in some fixed range, the agent gets a positive reward. If the CartPole is out of balance, the MDP transitions to an absorbing terminal state and gets 0 reward forever. Like in [Topin et al. \(2021\)](#), we use the gymnasium CartPole-v0 implementation [Towers et al. \(2024\)](#) of the CartPole MDP in which trajectories are truncated after 200 timesteps making the maximum cumulative reward, i.e. the optimal value of the RL objective when $\gamma = 1$, to be 200. The state features of the CartPole MDP are in $[-2, 2] \times [-2, 2] \times [-0.14, 0.14] \times [-1.4, 1.4]$.

IBMDP Authors define the associated IBMDP (definition 5) with $\zeta = -0.01$ and 4 information gathering actions. In addition to the original IBMDP paper, we also try $\zeta = 0.01$ and 3 information gathering actions. We use the same discount factor as the authors: $\gamma = 1$. We try two different approaches to limit the depth of decision tree policies to be at most 2: terminating trajectories if the agent takes too many information gathering actions in a row or simply giving a reward of -1 to the agent every time it takes an information gathering action past the depth limit. In practice, we could have tried an action masking approach, i.e. having a state dependent-action set, but we want to abide to the MDP formalism in order to properly understand direct interpretable approaches. We will also try IBMDPs where we do not limit the maximum depth for completeness.

Table 3: (Modified) PPO trained on 4×10^6 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()

Table 4: Top 5 hyperparameter configurations for modified DQN + IBMDP, bold font represent the original paper hyperparameters.

Rank	q	Depth control	Activation	Exploration	ζ	Mean 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

In tables 4 and 5 we report the top-5 hyperparameters for Modified RL baselines when learning partially observable IBMDP policies in terms of extracted decision tree policies performances in the CartPole MDP.

9.2 Experimental setup

Modified DQN and Modified PPO as mentioned above, the authors use the modified version of DQN from algorithm 2. We use the exact same hyperparameters for modified DQN as the authors when possible. We use the same layers width (128) and number of hidden layers (2), the same exploration strategy (ϵ -greedy with linearly decreasing value ϵ 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 ϵ value to 0.9). We use the same optimizer (RMSprop with hyperparameter 0.95 and learning rate 2.5×10^{-4}) to update the Q -networks. Authors did not share which DQN implementation they used so we use the stable-baselines3 one Raffin et al. (2021)³. Authors did not share which activation functions they used so we try both tanh and relu. For the modified PPO algorithm (algorithm 3), we can exactly match the authors hyperparameters since they

³We are cleaning our source code and will open source it as soon as possible

Table 5: Top 5 hyperparameter configurations for modified PPO + IBMDP, bold font represent the original paper hyperparameters.

Rank	q	Depth Control	Activation	ζ	Mean Final Performance
1	3	reward	relu	0.01	139
2	3	termination	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

Data: IBMDP $\mathcal{M}_{IB}\langle S \times O, A \cup A_{info}, (R, \zeta), (T, T_0, T_{info}) \rangle$, learning rate α , exploration rate ϵ , partially observable Q-network parameters θ , Q-network parameters ϕ , replay buffer \mathcal{B} , update frequency C

Result: deterministic memoryless policy π_{po}

Initialize partially observable Q-network parameters θ

Initialize Q-network parameters ϕ and target network parameters $\phi^- = \phi$

Initialize replay buffer $\mathcal{B} = \emptyset$

for each episode do

 Initialize downstream state features $\mathbf{s}_0 \sim T_0$

 Initialize observation $\mathbf{o}_0 = (L_1, U_1, \dots, L_p, U_p)$

for each step t do

 Choose action a_t using ϵ -greedy: $a_t = \operatorname{argmax}_a Q_\theta(\mathbf{o}_t, a)$ with prob. $1 - \epsilon$

 Take action a_t , observe r_t

 Store transition $(\mathbf{s}_t, \mathbf{o}_t, a_t, r_t, \mathbf{s}_{t+1})$ in \mathcal{B}

 Sample random batch $(\mathbf{s}_i, \mathbf{o}_i, a_i, r_i, \mathbf{s}_{i+1}) \sim \mathcal{B}$

$a' = \operatorname{argmax}_a Q_\theta(\mathbf{o}_i, a)$

$y_i = r_i + \gamma Q_\phi(\mathbf{s}_{i+1}, a')$ // Compute target $\phi \leftarrow \phi - \alpha \nabla_\phi (Q_\phi(\mathbf{s}_i, a_i) - y_i)^2$ // Update Q-network

$\theta \leftarrow \theta - \alpha \nabla_\theta (Q_\theta(\mathbf{o}_i, a_i) - y_i)^2$ // Update partially observable Q-network

if $t \bmod C = 0$ then

$\theta^- \leftarrow \theta$ // Update target network

end

$\mathbf{s}_t \leftarrow \mathbf{s}_{t+1}$

$\mathbf{o}_t \leftarrow \mathbf{o}_{t+1}$

end

end

$\pi_{po}(\mathbf{o}) = \operatorname{argmax}_a Q_\theta(\mathbf{o}, a)$ // Extract greedy policy

Algorithm 2: Modified Deep Q-Network. We highlight in green the changes to the standard DQN Mnih et al. (2015).

Data: IBMDP $\mathcal{M}_{IB}\langle S \times O, A \cup A_{info}, (R, \zeta), (T, T_0, T_{info}) \rangle$, learning rate α , policy parameters θ , clipping parameter ϵ , value function parameters ϕ

Result: Memoryless stochastic policy $\pi_{po\theta}$

Initialize policy parameters θ and value function parameters ϕ

for each episode do

 Generate trajectory $\tau = (\mathbf{s}_0, \mathbf{o}_0, a_0, r_0, \mathbf{s}_1, \mathbf{o}_1, a_1, r_1, \dots)$ following π_θ

for each timestep t in trajectory do

$G_t \leftarrow \sum_{k=t}^T \gamma^{k-t} r_k$ // Compute return $A_t \leftarrow G_t - V_\phi(\mathbf{s}_t)$ // Compute advantage

$r_t(\theta) \leftarrow \frac{\pi_{po\theta}(a_t|\mathbf{o}_t)}{\pi_{po\theta old}(a_t|\mathbf{o}_t)}$ // Compute probability ratio

$L_t^{CLIP} \leftarrow \min(r_t(\theta)A_t, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon)A_t)$ // Clipped objective

$\theta \leftarrow \theta + \alpha \nabla_\theta L_t^{CLIP}$ // Policy update $\phi \leftarrow \phi + \alpha \nabla_\phi (G_t - V_\phi(\mathbf{s}_t))^2$ // Value function update

end

$\theta_{old} \leftarrow \theta$ // Update old policy

end

Algorithm 3: Modified Proximal Policy Optimization

use the open source stable-baselines3 implementation of PPO. We match training budgets: we train modified DQN on 1 million timesteps and modified PPO on 4 million timesteps.

DQN and PPO We also benchmark the standard DQN and PPO when learning standard Markovian IBMDP policies $\pi : S \times O \rightarrow A \cup A_{info}$ and when learning standard $\pi : S \rightarrow A$ policies directly in the CartPole MDP. We summarize hyperparameters for the IBMDP and for the learning algorithms in appendices 1, 2 and 3.

Indirect methods We also compare modified RL algorithm to imitation learning. To do so, we use VIPER or Dagger to imitate greedy neural network policies obtained with standard DQN learning directly on CartPole. We use Dagger to imitate neural network policies obtained with the standard PPO learning directly on CartPole. For each indirect method, we imitate the neural network experts by fitting decision trees on 10000 expert transitions using the CART Breiman et al. (1984) implementation from scikit-learn Pedregosa et al. (2011) with default hyperparameters and maximum depth of 2 like in Topin et al. (2021).

9.2.1 Metrics

The key metric of this section is performance when controlling the CartPole, i.e, the average *undiscounted* cumulative reward of a policy on 100 trajectories (RL objective with $\gamma = 1$). For modified RL algorithms that learn a memoryless policy (or Q -function) in an IBMDP, we periodically extract the policy (or Q -function) and use algorithm 1 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 undiscounted cumulative reward. For RL applied to IBMDPs, since we can’t deploy learned policies directly to the downstream MDP as the state dimensions mismatch—such policies are $S \times O \rightarrow A \cup A_{info}$ but the MDP states are in S —we periodically evaluate those IBMDP policies in a copy of the IBMDP in which we fix $\zeta = 0$ ensuring that the copied IBMDP undiscounted cumulative rewards only account rewards from the CartPole MDP (non-zero rewards in the IBMDP only occur when a reward from the downstream MDP is given, i.e. when $a_t \in A$ in the IBMDP (definition 5)). Similarly, we do 100 trajectories of the extracted policies in the copied IBMDP and report the average undiscounted cumulative reward. For RL applied directly to the downstream MDP we can just periodically extract the learned policies and evaluate them on 100 CartPole trajectories.

Since imitation learning baselines train offline, i.e, on a fixed dataset, their performances cannot directly be reported on the same axis as RL baselines. For that reason, during the training of a standard RL baseline, we periodically extract the trained neural policy/ Q -function that we consider as the expert to imitate. Those experts are then imitated with VIPER or Dagger using 10 000 newly generated transitions and then fitted decision tree policies are then evaluated on 100 CartPole trajectories. We do not report the imitation learning objective values during VIPER or Dagger training. Every single combination of IBMDP and Modified RL hyperparameters is run 20 times. For standard RL on either an IBMDP or an MDP, we use the paper original hyperparameters when they were specified, with depth control using negative rewards, $\tanh()$ activations. We use 20 individual random seeds for every experiment in this chapter. Next, we present our results when reproducing Topin et al. (2021).

9.3 Results

9.3.1 How well do modified deep RL baselines learn in IBMDPs?

On figure 9a, we observe that modified DQN can learn in IBMDPs—the curves have an increasing trend—but we also observe that modified DQN finds poor decision tree policies for the CartPole MDP in average—the curves flatten at the end of the x-axis and have low y-values. In particular, the highest final y-value, among all the learning curves that could possibly correspond to the original paper modified DQN, correspond to poor performances on the CartPole MDP. On figure 9b, we observe that modified PPO finds decision tree policies with almost 150 cumulative rewards towards the end

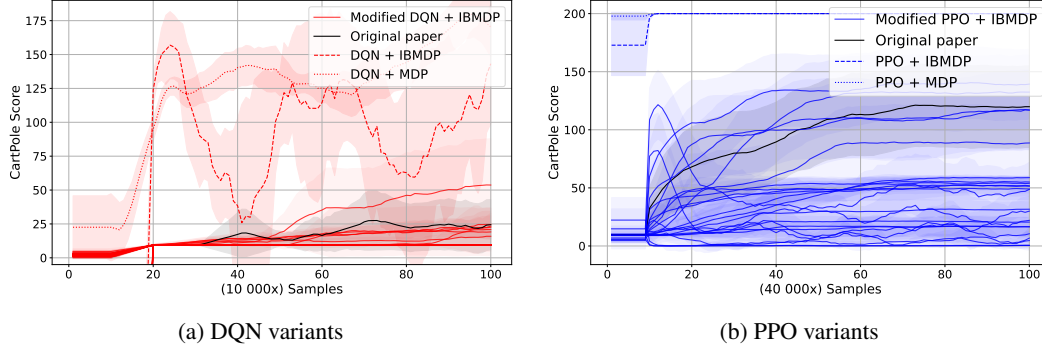


Figure 9: Comparison of modified reinforcement learning algorithms on different CartPole IBMDPs. (a) Shows variations of modified DQN and DQN (table 2), while (b) shows variations of modified PPO and PPO (table 3). For both algorithms, we give different line-styles for the learning curves when applied directly on the CartPole MDP versus when applied on the IBMDP to learn standard Markovian policies. We color the modified RL algorithm variant from the original paper in black. Shaded areas represent the confidence interval at 95% at each measure on the y-axis.

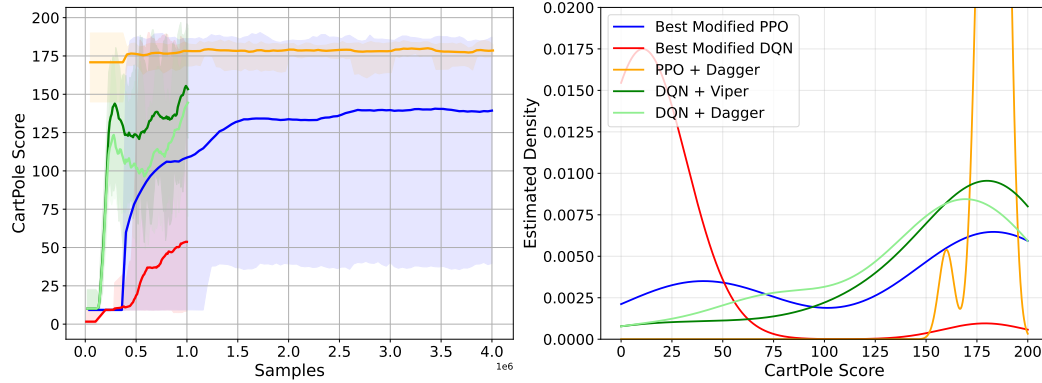


Figure 10: (left) Mean performance of the best-w.r.t. the RL objective for CartPole-modified RL + IBMDP combination. Shaded areas represent the min and max performance over the 20 seeds during training. (right) Corresponding score distribution of the final decision tree policies w.r.t. the RL objective for CartPole.

of training. The performance difference with modified DQN could be because we trained modified PPO longer, like in the original paper. However it could also be because DQN-like algorithms with those hyperparameters struggle to learn in CartPole (IB)MDPs. Indeed, we notice that for DQN-like baselines, learning seems difficult in general independently of the setting. On figures 9a and 9b, we observe that standard RL baselines (RL + IBMDP and RL + MDP), learn better CartPole policies in average than their modified counterparts that learn memoryless policies. On figure 9b, it is clear that for the standard PPO baselines, learning is super efficient and algorithms learn optimal policies with reward 200 in few thousands steps.

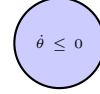
9.3.2 Which decision tree policies does direct reinforcement learning return for the CartPole MDP?

On figure 10, we isolate the best performing algorithms instantiations that learn decision tree policies for the CartPole MDP. We compare the best modified DQN and modified PPO to imitation learning baselines that use the surrogate imitation objective to find CartPole decision tree policies. We find that despite having poor performances in *average*, the modified deep reinforcement learning

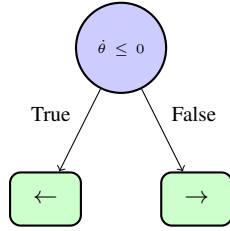
Most frequent modified PPO tree (12)



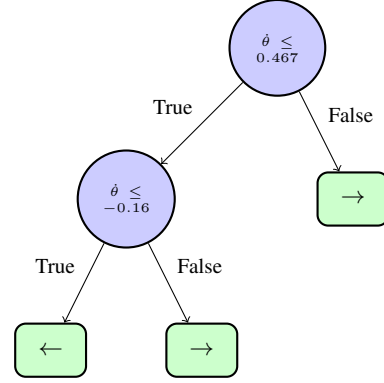
Most frequent modified DQN tree (9.5)



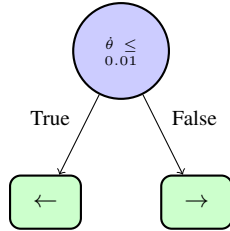
Best modified PPO tree (175)



Best modified DQN tree (160)



Typical imitated tree (185)



Best DQN + VIPER tree (200)

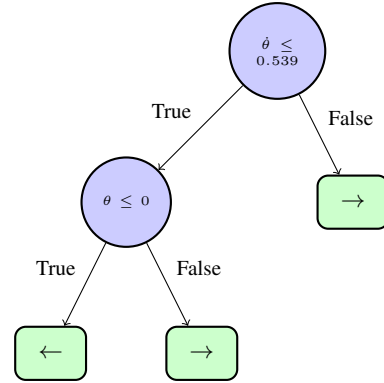


Figure 11: Trees obtained by modified deep RL in IBMDPs against trees obtained with imitation (RL objective value). θ and $\dot{\theta}$ are respectively the angle and the angular velocity of the pole

baselines can find very good decision tree policies as shown by the min-max shaded areas on the left of figure 10 and the corresponding estimated density of learned trees performances. However this is not desirable, a user typically wants an algorithm that can consistently find good decision tree policies. As shown by the estimated densities, indirect methods consistently find good decision tree policies (the higher modes of distributions are on the right of the plot). On the other hand, the decision tree policies returned by direct RL methods seem equally distributed on both extremes of the scores.

On figure 11, we present the best decision tree policies for CartPole returned by modified DQN and modified PPO. We used algorithm 1 to extract 20 trees from the 20 memoryless policies returned by the modified deep reinforcement learning algorithms over the 20 training seeds. We then plot the best tree for each baseline. Those trees get an average RL objective of roughly 175. Similarly, we plot a representative tree for imitation learning baseline as well as a tree that is optimal for CartPole w.r.t. the RL objective obtained with VIPER. Unlike for direct methods, the trees returned by imitation learning are extremely similar across seeds. 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 trees across seeds returned by modified RL baselines are “trivial” decision tree policies that either repeat the same downstream action forever or repeat the same IGA (definition 5) forever.

10 RL objective values calculations

Optimal depth-1 decision tree policy π_{τ_1} has one root node that tests $x \leq 1$ (respectively $y \leq 1$) and two leaf nodes \rightarrow and \downarrow . To compute $V_{\tau_1}^{\pi}(o_0)$, we compute the values of π_{τ_1} in each of the possible starting states $(s_0, o_0), (s_1, o_0), (s_2, o_0), (s_g, o_0)$ and compute the expectation over those. At initialization, when the downstream 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:

$$\begin{aligned} V^{\pi_{\tau_1}}(s_g, o_0) &= \zeta + \gamma + \gamma^2 \zeta + \gamma^3 \dots \\ &= \sum_{t=0}^{\infty} \gamma^{2t} \zeta + \sum_{t=0}^{\infty} \gamma^{2t+1} \\ &= \frac{\zeta + \gamma}{1 - \gamma^2} \end{aligned}$$

At initialization, in either of the downstream 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 :

$$\begin{aligned} V^{\pi_{\tau_1}}(s_0, o_0) &= \zeta + \gamma 0 + \gamma^2 V^{\pi_{\tau_1}}(s_g, o_0) \\ &= \zeta + \gamma^2 V^{\pi_{\tau_1}}(s_g, o_0) \\ &= V^{\pi_{\tau_1}}(s_2, o_0) \end{aligned}$$

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 :

$$\begin{aligned} V^{\pi_{\tau_1}}(s_1, o_0) &= \zeta + \gamma 0 + \gamma^2 V^{\pi_{\tau_1}}(s_2, o_0) \\ &= \zeta + \gamma^2 V^{\pi_{\tau_1}}(s_2, o_0) \\ &= \zeta + \gamma^2 (\zeta + \gamma^2 V^{\pi_{\tau_1}}(s_g, o_0)) \\ &= \zeta + \gamma^2 \zeta + \gamma^4 V^{\pi_{\tau_1}}(s_g, o_0) \end{aligned}$$

848 Since the probability of being in any downstream states at initialization given that the observation is
 849 \mathbf{o}_0 is simply the probability of being in any downstream states at initialization, we can write:

$$\begin{aligned} V^{\pi_{\tau_1}}(\mathbf{o}_0) &= \frac{1}{4}V^{\pi_{\tau_1}}(\mathbf{s}_g, \mathbf{o}_0) + \frac{2}{4}V^{\pi_{\tau_1}}(\mathbf{s}_2, \mathbf{o}_0) + \frac{1}{4}V^{\pi_{\tau_1}}(\mathbf{s}_1, \mathbf{o}_0) \\ &= \frac{1}{4}\frac{\zeta + \gamma}{1 - \gamma^2} + \frac{2}{4}(\zeta + \gamma^2\frac{\zeta + \gamma}{1 - \gamma^2}) + \frac{1}{4}(\zeta + \gamma^2\zeta + \gamma^4\frac{\zeta + \gamma}{1 - \gamma^2}) \\ &= \frac{1}{4}\frac{\zeta + \gamma}{1 - \gamma^2} + \frac{2}{4}(\frac{\zeta + \gamma^3}{1 - \gamma^2}) + \frac{1}{4}(\frac{\zeta + \gamma^5}{1 - \gamma^2}) \\ &= \frac{4\zeta + \gamma + 2\gamma^3 + \gamma^5}{4(1 - \gamma^2)} \end{aligned}$$

850 **Depth-0 decision tree:** has only one leaf node that takes a single downstream action indefinitely.
 851 For this type of tree the best reward achievable is to take actions that maximize the probability of
 852 reaching the objective \rightarrow or \downarrow . In that case the objective value of such tree is: In the goal state
 853 $G = (1, 0)$, the value of the depth-0 tree \mathcal{T}_0 is:

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

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

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

856 In the other states the policy never gets positive rewards; $V_{S_1}^{\mathcal{T}_0} = V_{S_2}^{\mathcal{T}_0} = 0$. Hence:

$$\begin{aligned} J(\mathcal{T}_0) &= \frac{1}{4}V_G^{\mathcal{T}_0} + \frac{1}{4}V_{S_0}^{\mathcal{T}_0} + \frac{1}{4}V_{S_1}^{\mathcal{T}_0} + \frac{1}{4}V_{S_2}^{\mathcal{T}_0} \\ &= \frac{1}{4}V_G^{\mathcal{T}_0} + \frac{1}{4}\gamma V_G^{\mathcal{T}_0} + 0 + 0 \\ &= \frac{1}{4}\frac{1}{1 - \gamma} + \frac{1}{4}\gamma\frac{1}{1 - \gamma} \\ &= \frac{1 + \gamma}{4(1 - \gamma)} \end{aligned}$$

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

$$V_{S_0}^{\mathcal{T}_u} = \zeta + \gamma\zeta + \gamma^2 0 + \gamma^3 V_G^{\mathcal{T}_1}$$

862

$$\begin{aligned} V_{S_1}^{\mathcal{T}_u} &= \zeta + \gamma\zeta + \gamma^2 0 + \gamma^3 V_{S_0}^{\mathcal{T}_u} \\ &= \zeta + \gamma\zeta + \gamma^2 0 + \gamma^3(\zeta + \gamma\zeta + \gamma^2 0 + \gamma^3 V_G^{\mathcal{T}_1}) \\ &= \zeta + \gamma\zeta + \gamma^3\zeta + \gamma^4\zeta + \gamma^6 V_G^{\mathcal{T}_1} \end{aligned}$$

863 We get:

$$\begin{aligned}
J(\mathcal{T}_u) &= \frac{1}{4}V_G^{\mathcal{T}_u} + \frac{1}{4}V_{S_0}^{\mathcal{T}_u} + \frac{1}{4}V_{S_1}^{\mathcal{T}_u} + \frac{1}{4}V_{S_2}^{\mathcal{T}_u} \\
&= \frac{1}{4}V_G^{\mathcal{T}_1} + \frac{1}{4}(\zeta + \gamma\zeta + \gamma^3V_G^{\mathcal{T}_1}) + \frac{1}{4}(\zeta + \gamma\zeta + \gamma^3\zeta + \gamma^4\zeta + \gamma^6V_G^{\mathcal{T}_1}) + \frac{1}{4}V_{S_2}^{\mathcal{T}_1} \\
&= \frac{1}{4}\left(\frac{\zeta + \gamma}{1 - \gamma^2}\right) + \frac{1}{4}\left(\frac{\gamma\zeta + \gamma^4 + \zeta - \gamma^2\zeta}{1 - \gamma^2}\right) + \frac{1}{4}(\zeta + \gamma\zeta + \gamma^3\zeta + \gamma^4\zeta + \gamma^6V_G^{\mathcal{T}_1}) + \frac{1}{4}V_{S_2}^{\mathcal{T}_1} \\
&= \frac{1}{4}\left(\frac{\zeta + \gamma}{1 - \gamma^2}\right) + \frac{1}{4}\left(\frac{\gamma\zeta + \gamma^4 + \zeta - \gamma^2\zeta}{1 - \gamma^2}\right) + \frac{1}{4}\left(\frac{\zeta + \gamma\zeta - \gamma^2\zeta - \gamma^5\zeta + \gamma^6\zeta + \gamma^7}{1 - \gamma^2}\right) + \frac{1}{4}V_{S_2}^{\mathcal{T}_1} \\
&= \frac{1}{4}\left(\frac{\zeta + \gamma}{1 - \gamma^2}\right) + \frac{1}{4}\left(\frac{\gamma\zeta + \gamma^4 + \zeta - \gamma^2\zeta}{1 - \gamma^2}\right) + \frac{1}{4}\left(\frac{\zeta + \gamma\zeta - \gamma^2\zeta - \gamma^5\zeta + \gamma^6\zeta + \gamma^7}{1 - \gamma^2}\right) + \frac{1}{4}\left(\frac{\zeta + \gamma^3}{1 - \gamma^2}\right) \\
&= \frac{\zeta(4 + 2\gamma - 2\gamma^2 - \gamma^5 + \gamma^6) + \gamma + \gamma^3 + \gamma^4 + \gamma^7}{4(1 - \gamma^2)}
\end{aligned}$$

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

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

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

$$\begin{aligned}
J(\mathcal{T}_2) &= \frac{1}{4}V_G^{\mathcal{T}_2} + \frac{2}{4}V_{S_2}^{\mathcal{T}_2} + \frac{1}{4}V_{S_1}^{\mathcal{T}_2} \\
&= \frac{1}{4}V_G^{\mathcal{T}_2} + \frac{2}{4}(\zeta + \gamma\zeta + \gamma^20 + \gamma^3V_G^{\mathcal{T}_2}) + \frac{1}{4}(\zeta + \gamma\zeta + \gamma^20 + \gamma^3\zeta + \gamma^4\zeta + \gamma^50 + \gamma^6V_G^{\mathcal{T}_2}) \\
&= \frac{\zeta(3 + 3\gamma) + \gamma^2 + \gamma^5 + \gamma^8}{4(1 - \gamma^3)}
\end{aligned}$$

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

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

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

874 Solving these equations:

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

875

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

876

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

877 11 Training curves

878 We plot the sub-optimality gaps during training, w.r.t. the RL objective (definition 2),
879 between the learned memoryless policy and the optimal deterministic memoryless policy:
880 $|\mathbb{E}[V^{\pi^*}(s_0, o_0)|s_0 \sim T_0] - \mathbb{E}[V^{\pi}(s_0, o_0)|s_0 \sim T_0]|$. Because we know the whole POIBMDP
881 model that we can represent exactly as tables, and because we know for each ζ the RL objective value
882 of the optimal deterministic memoryless policy (figure 2), we can report the *exact* sub-optimality
883 gaps. In figure 12, we plot the sub-optimality gaps—averaged over 100 seeds—of learned policies
884 during training. We do so for 200 different POIBMDPs where we change the reward for informa-
885 tion gathering actions: we sample 200 ζ values uniformly in $[-1, 2]$. In figure 12, a different color
886 represents a different POIBMDP.

887 Recall from figure 2 that for: (i) $\zeta \in [-1, 0]$, the optimal deterministic memoryless policy is a depth-
888 0 tree, (ii) $\zeta \in]0, 1[$, the optimal deterministic memoryless policy is a depth-1 tree, and (iii) $\zeta \in [1, 2]$,
889 the optimal deterministic memoryless policy is a “infinite” tree that contains infinite number of
890 internal nodes. We observe that, despite all sub-optimality gaps converging independently of the ζ
891 values, not all algorithms in all POIBMDPs fully minimize the sub-optimality gap. In particular, all
892 algorithms seem to consistently minimize the gap, i.e. learn the optimal policy or Q-function, only
893 for $\zeta \in [1, 2]$ (all the yellow lines go to 0). However, we are interested in the range $\zeta \in]0, 1[$ where
894 the optimal decision tree policy is non-trivial, i.e. not taking the same action forever. In that range,
895 no baseline consistently minimizes the sub-optimality gap.

896 12 Tabular RL algorithmic details for POIBMDPs

897 12.1 Training with the best hyperparameters

898 13 POIBMDPs for classification tasks

899 Let us show that, POIBMDPs associated with MDPs encoding supervised learning tasks, are in fact
900 MDPs themselves. Let us define such supervised learning MDPs in the context of a classification
901 task (this definition extends trivially to regression tasks).

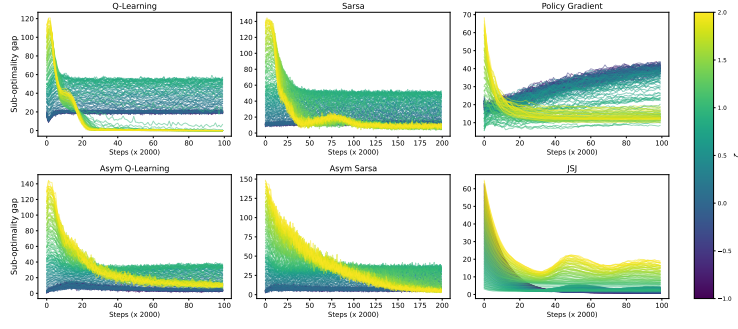


Figure 12: (Asymmetric) reinforcement learning in POIBMDPs. In each subplot, each single line is colored by the value of ζ in the corresponding POIBMDP in which learning occurs. Each single learning curve represent the sub-optimality gap averaged over 100 seeds.

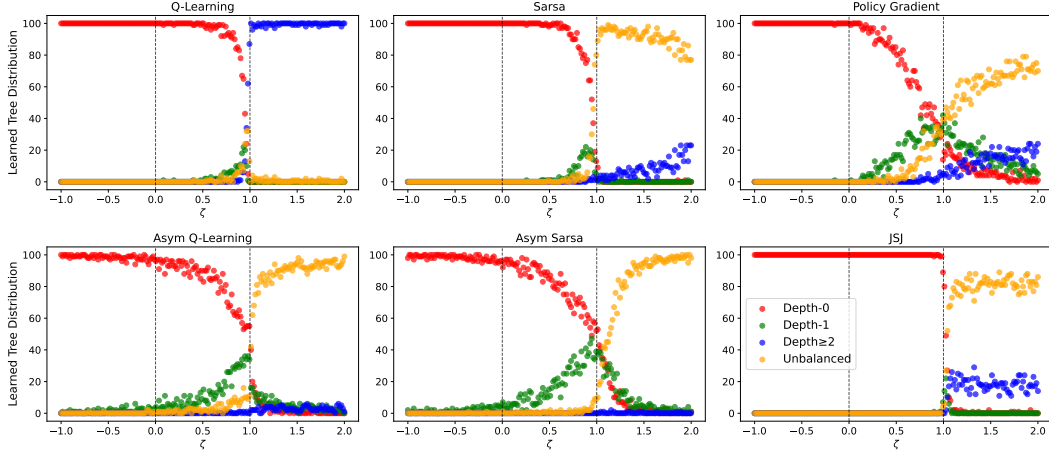
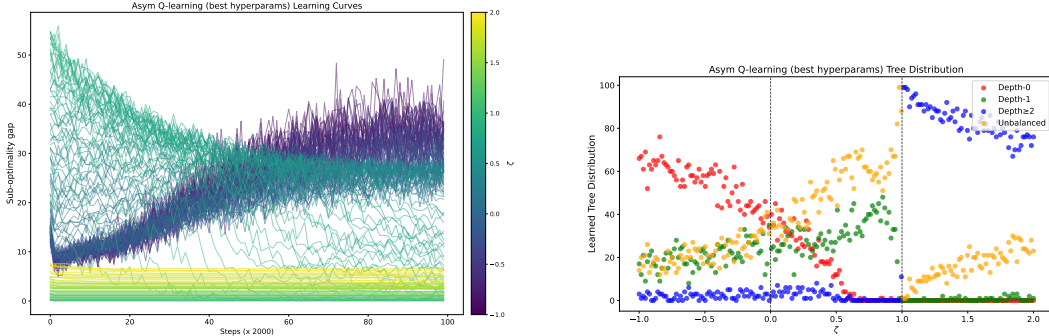


Figure 13: Distributions of final tree policies learned across the 100 seeds. For each ζ value, there are four colored points. Each point represent the share of depth-0 trees (red), depth-1 trees (green), unbalanced depth-2 trees (orange) and depth-2 trees (blue).



(a) Learning curves for asymmetric Q-learning with good hyperparameters.

(b) Trees distributions for asymmetric Q-learning with good hyperparameters

Figure 14: Analysis of the top-performing asymmetric Q-learning instantiation. (left) Learning curves, and (right) tree distributions across different POIBMDP configurations.

Algorithm 4: Asymmetric Q-Learning. We highlight in green the differences with the standard Q-learning [Watkins & Dayan \(1992\)](#)

Data: A POMDP, learning rates α_u, α_q , exploration prob. ϵ

Result: $\pi : O \rightarrow A$

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

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

for each episode do

 Initialize state $x_0 \sim T_0$

 Initialize observation $\mathbf{o}_0 \sim \Omega(x_0)$

for each step t do

 Choose action a_t using ϵ -greedy: $a_t = \operatorname{argmax}_a Q(\mathbf{o}_t, a)$ with prob. $1 - \epsilon$

 Take action a_t , observe $r_t = R(\mathbf{x}_t, a_t)$, $x_{t+1} \sim T(x_t, a_t)$, and $\mathbf{o}_{t+1} \sim \Omega(x_{t+1})$

$y \leftarrow r + \gamma U(\mathbf{x}_{t+1}, \operatorname{argmax}_{a'} Q(\mathbf{o}_{t+1}, a'))$

$U(\mathbf{x}_t, a_t) \leftarrow (1 - \alpha_u)U(\mathbf{x}_t, a_t) + \alpha_u y$

$Q(\mathbf{o}_t, a_t) \leftarrow (1 - \alpha_q)Q(\mathbf{o}_t, a_t) + \alpha_q y$

$x_t \leftarrow x_{t+1}$

$\mathbf{o}_t \leftarrow \mathbf{o}_{t+1}$

end

end

$\pi(o) = \operatorname{argmax}_a Q(\mathbf{o}, a)$

Algorithm 5: Asymmetric Sarsa

Data: POMDP $\mathcal{M}_{po} = \langle X, O, A, R, T, T_0, \Omega \rangle$, learning rates α_u, α_q , exploration rate ϵ

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 $\mathbf{o} \in O, a \in A$

for each episode do

 Initialize state $x_0 \sim T_0$

 Initialize observation $\mathbf{o}_0 \sim \Omega(x_0)$

 Choose action a_0 using ϵ -greedy: $a_0 = \operatorname{argmax}_a Q(\mathbf{o}_0, a)$ with prob. $1 - \epsilon$

for each step t do

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

 Choose action a_{t+1} using ϵ -greedy: $a_{t+1} = \operatorname{argmax}_a Q(\mathbf{o}_{t+1}, a)$ with prob. $1 - \epsilon$

$y \leftarrow r + \gamma U(x_{t+1}, a_{t+1})$ // TD target using actual next action

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

$Q(\mathbf{o}_t, a_t) \leftarrow (1 - \alpha_q)Q(\mathbf{o}_t, a_t) + \alpha_q y$

$x_t \leftarrow x_{t+1}$

$\mathbf{o}_t \leftarrow \mathbf{o}_{t+1}$

$a_t \leftarrow a_{t+1}$

end

end

$\pi(\mathbf{o}) = \operatorname{argmax}_a Q(\mathbf{o}, a)$ // Extract greedy policy

Algorithm 6: Asymmetric policy gradient algorithm. Uses Monte Carlo estimates of the average reward value functions to perform policy improvements.

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

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

Initialize policy parameters θ

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 π_θ

for each timestep t **in trajectory** τ_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 $a = \operatorname{argmax}_{a'} Q(o, a')$, 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

Table 6: Summary of RL baselines Hyperparameters

algorithm	Problem	Hyperparameters comb.
Policy Gradient	PO/IB/MDP	420
JSJ	POIBMDP	15
Q-learning	PO/IB/MDP	192
Asym Q-learning	POIBMDP	768
Sarsa	PO/IB/MDP	192
Asym Sarsa	POIBMDP	768

Table 7: 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

Table 8: 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

Table 9: 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

Table 10: 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

Table 11: 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

Table 12: 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

Table 13: 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

Table 14: 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 15: Asymmetric sarsa hyperparameters (768 combinations each run 10 times)

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	learning rate for the Q-function
Learning Rate Q	0.001, 0.005, 0.01, 0.1	learning rate for the partial observation dependent Q-function
Optimistic	True, False	Optimistic initialization

Hyperparameter	Asym Q-learning (10/10)	Asym Sarsa (10/10)	PG (4/10)
epsilon_start	1.0	1.0	-
epsilon_decay	0.99	0.99	-
batch_size	1	1	-
lambda_	0.0	0.0	-
lr_o	0.01	0.1	-
lr_v	0.1	0.005	-
optimistic	False	False	-
lr	-	-	0.05
tau	-	-	0.1
eps	-	-	0.1
n_steps	-	-	2000

Table 16: Best hyperparameters for each algorithm on the POIBMDP problem

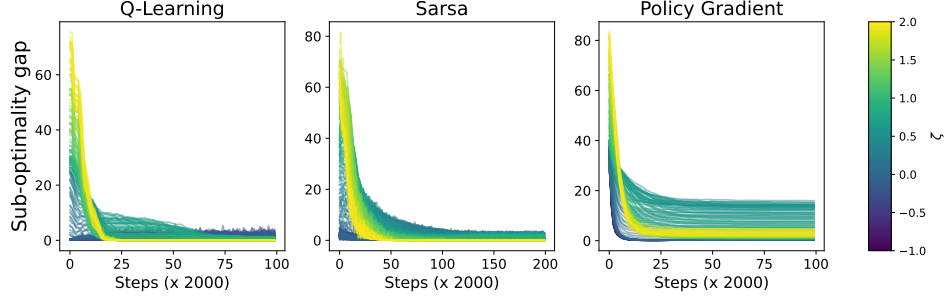


Figure 15: We reproduce the same plot as in figure 12 for classification POIBMDPs. Each individual curve is the sub-optimality gap of the learned policy during training averaged over 100 runs for a single ζ value.

Definition 9 (Classification Markov decision process). Given a set of N examples denoted $\mathcal{E} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$ where each datum $\mathbf{x}_i \in \mathcal{X}$ is described by a set of p features x_{ij} with $1 \leq j \leq p$, and $y_i \in \mathbb{Z}^m$ is the label associated with \mathbf{x}_i , a classification Markov decision Process is an MDP $\langle S, A, R, T, T_0 \rangle$ (definition 1). The state space is $S = \{\mathbf{x}_i\}_{i=1}^N$, the set of training 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 $R(s = \mathbf{x}_i, a) = 1_{\{a=y_i\}}$. 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 $T_0(s_0 = s) = \frac{1}{N}$.

One can be convinced that policies that maximize the RL objective (definition 2) in classification MDPs are classifiers that maximize the prediction accuracy because $\sum_{i=1}^N 1_{\pi(\mathbf{x}_i)=y_i} = \sum_{i=1}^N R(\mathbf{x}_i, \pi(\mathbf{x}_i))$. We defer the formal proof in the next part of the manuscript in which we extensively study supervised learning problems.

Now let us show that associated POIBMDPs are in fact MDPs. We show this by construction.

Definition 10 (Classification POIBMDP). Given a classification MDP $\langle \{\mathbf{x}_i\}_{i=1}^N, \mathbb{Z}^m, R, T, T_0 \rangle$ (definition 9), and an associated POIBMDP $\langle S, O, A, A_{info}, R, \zeta, T_{info}, T, T_0 \rangle$ (definition 7), a classification POIBMDP is an MDP (definition 1):

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

O is the set of possible observations in $[L_1, U_1] \times \dots \times [L_p, U_p] \times [L_1, U_1] \times \dots \times [L_p, U_p]$ where L_j is the minimum value of feature j over all data \mathbf{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 refinements in A_{info} . The reward for assigning label $a \in \mathbb{Z}^m$ when observing some observation $\mathbf{o} = (L'_1, U'_1, \dots, L'_p, U'_p)$ is the proportion of training data satisfying the bounds and having label a : $R(\mathbf{o}, a) = \frac{|\{\mathbf{x}_i: L'_j \leq x_{ij} \leq U'_j \forall i, j\} \cap \{\mathbf{x}_i: y_i = a \forall i\}|}{|\{\mathbf{x}_i: L'_j \leq x_{ij} \leq U'_j \forall i, j\}|}$. The reward for taking an information gathering action that refines bounds is ζ . The transition function is $\mathcal{P} : O \times (\mathbb{Z}^m \cup A_{info}) \rightarrow \Delta(O)$. When $a \in \mathbb{Z}^m$, $\mathcal{P}(\mathbf{o}, a, (L_1, U_1, \dots, L_p, U_p)) = 1$ (reset to full bounds). When $a = (k, v) \in A_{info}$, from $\mathbf{o} = (L'_1, U'_1, \dots, L'_p, U'_p)$, the MDP will transit to $\mathbf{o}_{left} = (L'_1, U'_1, \dots, L'_k, v, \dots, L'_p, U'_p)$ (resp. $\mathbf{o}_{right} = (L'_1, U'_1, \dots, U'_k, v, \dots, L'_p, U'_p)$) with probability $\frac{|\{\mathbf{x}_i: L'_j \leq x_{ij} \leq U'_j \forall j \wedge x_{ik} \leq v\}|}{|\{\mathbf{x}_i: L'_j \leq x_{ij} \leq U'_j \forall j\}|}$ (resp. $\frac{|\{\mathbf{x}_i: L'_j \leq x_{ij} \leq U'_j \forall j \wedge x_{ik} > v\}|}{|\{\mathbf{x}_i: L'_j \leq x_{ij} \leq U'_j \forall j\}|}$).

Those classification POIBMDPs are essentially MDPs with stochastic transitions. It means that deterministic memoryless policies $O \rightarrow A \cup A_{info}$ are in fact Markovian policy for those classification POIBMDPs. More importantly, it means that, for a given γ and ζ , if we were to know the whole POIBMDP model, we could use planning, to compute *optimal* decision tree policies.