

Disclaimer/Publisher's Note: The statements, opinions, and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions, or products referred to in the content.

Article

# Iterative Oblique Decision Trees Deliver Explainable RL Models

Raphael C. Engelhardt <sup>1</sup>, Marc Oedingen <sup>1</sup>, Moritz Lange <sup>2</sup>, Laurenz Wiskott <sup>2</sup> and Wolfgang Konen <sup>1,\*</sup>

<sup>1</sup> Cologne Institute of Computer Science, Faculty of Computer Science and Engineering Science, TH Köln, Gummersbach, Germany; {Raphael.Engelhardt,Marc.Oedingen,Wolfgang.Konen}@th-koeln.de  
<sup>2</sup> Institute for Neural Computation, Faculty of Computer Science, Ruhr-University Bochum, Bochum, Germany; {Moritz.Lange,Laurenz.Wiskott}@ini.rub.de  
\* Correspondence: Wolfgang.Konen@th-koeln.de

**Abstract:** The demand for explainable and transparent models increases with the continued success of reinforcement learning. In this article, we explore the potential of generating shallow decision trees (DT) as simple and transparent surrogate models for opaque deep reinforcement learning (DRL) agents. We investigate three algorithms for generating training data for axis-parallel and oblique DTs with the help of DRL agents (“oracles”) and evaluate these methods on classic control problems from OpenAI Gym. The results show that one of our newly developed algorithms, the iterative training, outperforms traditional sampling algorithms, resulting in well-performing DTs that often even surpass the oracle from which they were trained. Even higher dimensional problems can be solved with surprisingly shallow DTs. We discuss the advantages and disadvantages of different sampling methods and insights into the decision-making process made possible by the transparent nature of DTs. Our work contributes to the development of not only powerful but also explainable RL agents and highlights the potential of DTs as a simple and effective alternative to complex DRL models.

**Keywords:** Reinforcement learning, Decision tree, Explainable AI, Rule learning

## 1. Introduction

One of the most significant drawbacks of powerful deep reinforcement learning (DRL) algorithms is their opacity. The well-performing decision-making process is buried in the depth of artificial neural networks, which might constitute a major barrier to the application of reinforcement learning (RL) in various areas.

In a recent publication [1], we proposed an algorithm for obtaining simple decision trees (DT) from trained DRL agents (“oracles”). This approach has notable advantages:

1. It is conceptually simple, as it translates the problem of explainable RL into a supervised learning setting.
2. DTs are fully transparent and (at least for limited depth) offer a set of easily understandable rules.
3. The approach is oracle-agnostic: it does not rely on the agent being trained by a specific RL algorithm, as only a training set of state-action pairs is required.

Although the algorithm in [1] was successful on some problems, it failed on others because no well-performing shallow DT could be found. In this work, we investigate the reasons for those failures and propose three algorithms to generate training samples for the DTs: The first one is entirely based on evaluation episodes of the DRL agent (described in [1] and evaluated further in this article), the second one applies random sampling on a subregion of the observation space, and the third approach is an iterative algorithm involving the DRL agent’s predictions for regions explored by the DT. We test all approaches on seven environments of varying dimensionality, including all classic control problems from OpenAI Gym. Our results show that DTs of shallow depth can solve all environments, that DTs can even surpass DRL agents, and that the third algorithm (the iterative approach) outperforms

the other two methods in finding simple and well-performing trees. As a particularly surprising result, we show that an 8-dimensional control problem (LUNARLANDER) can be successfully solved with a DT of only depth 2.

We discuss the advantages and disadvantages of the three algorithms’ computational complexity and how DTs’ transparent, simple nature provides interesting options for explainability that are not immediately applicable to opaque DRL models.

The remainder of the article is structured as follows: In Section 2, we present related work. Section 3 contains a detailed description of our methods and the experimental setup. In Section 4, we present our results and discuss the explainability of DTs. We place our results in a broader context and discuss future research in Section 5 before drawing a short conclusion in Section 6.

2. Related Work

In recent years, interest in explainable and interpretable algorithms for Deep Learning algorithms has grown steadily. This is shown by review articles for explainable AI (XAI) and explainable machine learning in general [2–4] and, more recently, by review articles for explainable Reinforcement Learning (XRL) in particular [5–7]. Lundberg et al. [8] give an overview of efficient algorithms for generating explanations, especially for trees.

A variety of rule deduction methods have been developed over the years. Liu et al. [9] apply a mimic learning approach to RL using comparably complex Linear Model U-trees to approximate a Q-function. Our approach differs insofar as the DTs we obtain are arguably simpler and represent a transparent policy, translating state into action. Mania et al. [10] propose Augmented Random Search, another algorithm for finding linear models that solve RL problems. Coppens et al. [11] distill PPO agents’ policy networks into Soft Decision Trees [12] to get insights. Another interesting approach to explainable RL, which does not use DTs, is the approach by Verma et al. [13] called Programmatically Interpretable RL (PIRL) through Neurally Directed Program Synthesis (NDPS). In this method, the DRL guides the search of a policy consisting of specific operators and input variables to obtain interpretable rules mimicking PID controllers. As Verma et al. [13] also test their algorithms on OpenAI Gym’s classic control problems, a direct comparison is possible and was made in [1]. However, their results were not reproducible by us, and we show how our algorithms yield better results at low DT depths than the reported ones. Zilke et al. [14] propose an algorithm for distilling DTs from neural networks. With their method “DeepRED”, the authors translate each layer of a feedforward network into rules and aim to simplify those. Qiu and Zhu [15] present an algorithm that allows policy architectures to be learned and that does not require a trained oracle.

The work presented here extends our earlier research in Engelhardt et al. [1]. We apply it to additional environments and introduce new sampling algorithms to overcome the drawbacks of [1]. To the best of our knowledge, no other algorithms are available in the current state of the art to construct that simple DTs from trained DRL agents for a large variety of RL environments.

Concerning the iterative learning of trees outside of RL, well-known iterative meta-heuristics like boosting and its particular form AdaBoost [16,17] exist, which are often applied to trees. However, boosting is iterative with respect to weak learners and outputs an ensemble of them, while our approach is iterative in the data, retrains many trees, and outputs only the best-performing one.

Rudin [18] has emphasized the importance of building inherently interpretable models instead of trying to explain black box models. While Rudin [18] applied this to standard machine learning models for application areas like criminal justice, healthcare, and computer vision, we try to take her valuable arguments over to the area of control problems and RL and show how we can substitute black box models with inherently transparent DTs without loss of performance. Although we use a DRL model to construct our DT, the final DT can be used in operation and interpreted without referencing the DRL model.

### 3. Methods

Our experiments were implemented in Python programming language version 3.9.16 using a variety of different packages and methods described in the following<sup>1</sup>.

#### 3.1. Environments

We test our different approaches of training DTs from DRL agents on all classic control problems and the LUNARLANDER challenge offered by OpenAI Gym [19], and on the CARPOLE-SWINGUP problem as implemented in [20]. This selection of environments offers a variety of simple but not trivial control tasks with continuous, multidimensional tabular observables and one-dimensional discrete or continuous actions. Table 1 gives an overview of the different environments used in this article. For the two environments that do not have predefined thresholds at which they are considered solved, we manually assigned one corresponding to visually convincing episodes of quick and stable swing-ups.

**Table 1.** Overview of environments used. A task is considered solved when the average return in 100 evaluation episodes is  $\geq R_{solved}$ . Environments with undefined official threshold  $R_{solved}$  have been assigned a reasonable one. This number is then reported in parenthesis.

Environment	Observation space ( $\mathcal{O}$ )	Action space ( $\mathcal{A}$ )	$R_{solved}$
Acrobot-v1	First angle $\cos(\theta_1) \in [-1, 1]$ , First angle $\sin(\theta_1) \in [-1, 1]$ , Second angle $\cos(\theta_2) \in [-1, 1]$ , Second angle $\sin(\theta_2) \in [-1, 1]$ , Angular velocity $\dot{\theta}_1 = \omega_1 \in [-4\pi, 4\pi]$ , Angular velocity $\dot{\theta}_2 = \omega_2 \in [-9\pi, 9\pi]$	$a \in \text{Apply torque}\{-1 (0), 0 (1), 1 (2)\}$	-100
CartPole-v1	Position cart $x \in [-4.8, 4.8]$ , Velocity cart $v \in (-\infty, \infty)$ , Angle pole $\theta \in [-0.42, 0.42]$ , Velocity pole $\omega \in (-\infty, \infty)$	$a \in \text{Accelerate}\{\text{left} (0), \text{right} (1)\}$	475
CartPole-SwingUp-v1	Position cart $x \in [-2.4, 2.4]$ , Velocity cart $v \in (-\infty, \infty)$ , Angle pole $\cos(\theta) \in [-1, 1]$ , Angle pole $\sin(\theta) \in [-1, 1]$ , Angular velocity pole $\omega \in (-\infty, \infty)$	$a \in \text{Accelerate}\{\text{left} (0), \text{not} (1), \text{right} (2)\}$	undef. (840)
LunarLander-v2	Position $x \in [-1.5, 1.5]$ , Position $y \in [-1.5, 1.5]$ , Velocity $v_x \in [-5, 5]$ , Velocity $v_y \in [-5, 5]$ , Angle $\theta \in [-\pi, \pi]$ , Angular velocity $\dot{\theta} = \omega \in [-5, 5]$ , Contact left leg $l_l \in \{0, 1\}$ , Contact right leg $l_r \in \{0, 1\}$	$a \in \text{Fire}\{\text{not} (0), \text{left engine} (1), \text{main engine} (2), \text{right engine} (3)\}$	200
MountainCar-v0	Position $x \in [-1.2, 0.6]$ , Velocity $v \in [-0.07, 0.07]$	$a \in \text{Accelerate}\{\text{left} (0), \text{not} (1), \text{right} (2)\}$	-110
MountainCar Continuous-v0	Same as MountainCar-v0	$a \in \text{Accelerate} [-1, 1]$	90
Pendulum-v1	Angle $\cos(\theta) \in [-1, 1]$ , Angle $\sin(\theta) \in [-1, 1]$ , Angular velocity $\omega \in [-8, 8]$	$a \in \text{Apply torque} [-2, 2]$	undef. (-170)

<sup>1</sup> The experiments can be found in the Github repository <https://github.com/MarcOedingen/Iterative-Oblique-Decision-Trees>

### 3.2. Deep Reinforcement Learning

Our algorithms require a DRL agent to predict actions for given observations. We train DRL agents (“oracles”) using PPO [21], DQN [22], and TD3 [23] algorithms as implemented in the Python DRL framework Stable-Baselines3 (SB3) [24]. It is necessary to test all those algorithms because no single DRL algorithm solves all environments used in this paper. For each environment, we select one of the successful DRL agents (see Table 2a).

### 3.3. Decision Trees

For the training of DTs, we rely on two algorithms representing two “families” of DTs: The widely-used Classification and Regression Trees (CART) as described by Breiman *et al.* [25] and implemented in [26] partition the observation space with axis-parallel splits. The decision rules use only one observable at a time, which makes them particularly easy to interpret. On the other hand, Oblique Predictive Clustering Trees (OPCT), as described and implemented in [27], can generate DTs where each decision rule compares the weighted sum of the features  $o_i$  to a threshold  $\tau$ :

$$w_1 o_1 + w_2 o_2 + \dots + w_n o_n \leq \tau \quad (1)$$

Such rules subdivide the feature space with oblique splits (tilted hyperplanes), allowing to describe more complex partitions with fewer rules and, consequently, lower depth. However, the oblique split rules make the interpretation of such trees more challenging. We will address this topic in Section 4.4.

### 3.4. DT Training Methods

In this subsection, we describe our different algorithms to generate datasets of samples for the training of DTs.

#### 3.4.1. Episode Samples (EPS)

The basic approach has been described in detail in [1]. In brief, it consists of three steps:

1. A DRL agent (“oracle”) is trained to solve the problem posed by the studied environment.
2. The oracle acting according to its trained policy is evaluated for a set number of episodes. At each time step, the state of the environment and action of the agent are logged until a total of  $n_s$  samples are collected.
3. A decision tree (CART or OPCT) is trained from the samples collected in the previous step.

#### 3.4.2. Bounding Box (BB)

Detailed investigations of the failures of DTs trained on episode samples for the Pendulum environment show that a well-performing oracle only visits a very restricted region of the observation space and oversamples specific subregions: A good oracle is able to swing the pendulum in the upright position quickly and keeps the system in a small region around the corresponding point in the state space for the rest of the episode’s duration.

To counteract this issue, we investigated an approach based on querying the oracle at random points in the observation space. The algorithm consists of three steps:

1. The oracle is evaluated for a certain number of episodes. Let  $L_i$  and  $U_i$  be the lower and upper bound of the visited points in the  $i^{\text{th}}$  dimension of the observation space and  $\text{IQR}_i$  their interquartile range (difference between the 75<sup>th</sup> and the 25<sup>th</sup> percentile).
2. Based on the statistics of visited points in the observation space from the previous step, we take a certain number  $n_s$  of samples from the uniform distribution within a hyper-rectangle of side lengths  $[L_i - 1.5 \cdot \text{IQR}_i, U_i + 1.5 \cdot \text{IQR}_i]$ . The side length is clipped, should it exceed the predefined boundaries of the environment.

3. For all investigated depths  $d$ , OPCTs are trained from the dataset consisting of the  $n_s$  samples and the corresponding actions predicted by the oracle.

It should be noted that this algorithm has its peculiarities: Randomly chosen points in the observation space might correspond to inconsistent states. While this is evident for the example of PENDULUM (the observations contain  $\cos \theta$  and  $\sin \theta$  (see Table 1), which was solved here by sampling  $\theta$  and generating observations from that), less apparent dependencies between observables may occur in other cases.

### 3.4.3. Iterative Training of Explainable RL Models (ITER)

The DTs in [1] were trained solely based on samples from the oracle episodes. In the iterative algorithm, ITER, the tree with the best performance (from all iterations so far) is used as a generator for new observations.

The procedure is shown in Algorithm 1 and can be described as follows: The initial tree is trained on samples of oracle evaluation episodes. We use these initial samples to fit the tree to the oracle's predictions and use the tree's current, somewhat imperfect decision-making to generate observations in areas the oracle does maybe not access anymore. Often, the points in the observation space that are essential for the formation of a successful tree lie a bit beside oracle trajectories. Generating samples next to these trajectories with a not yet fully trained tree allows the next iteration's tree to fit its decisions to these previously unknown areas.

---

#### Algorithm 1 Iterative Training of Explainable RL Models (ITER)

---

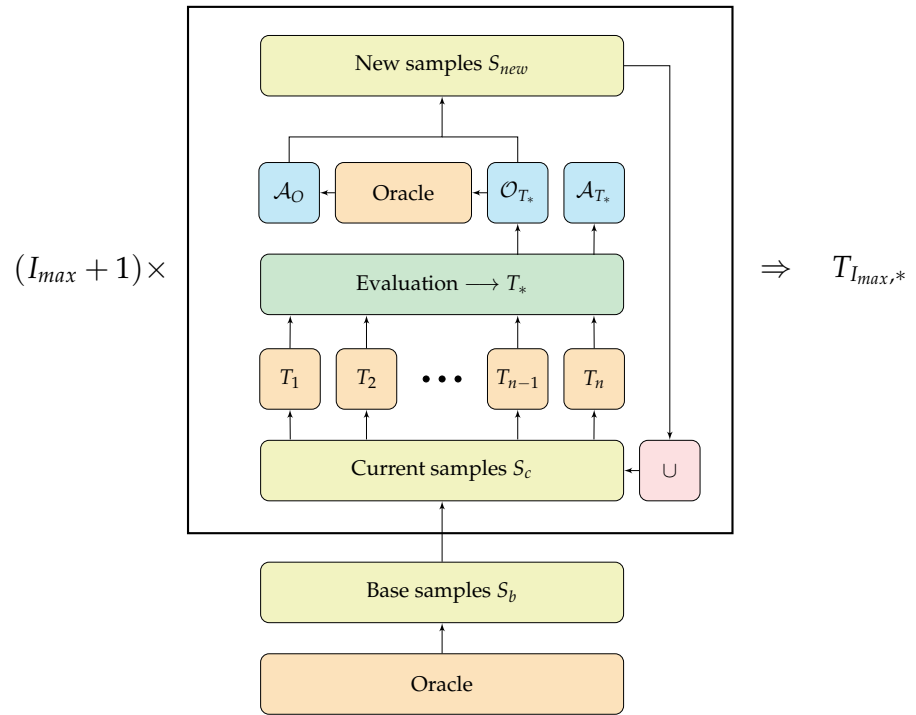
**Require:** Oracle policy  $\pi_O$ , maximum iteration  $I_{max}$ , number  $n_b$  of base samples, number of trees  $T_{max}$ , evaluation episodes  $n_{eps}$ , DT depth  $d$ , number of samples added in each iteration  $n_{samp}$ .  $\mathcal{O}$ : observations,  $\mathcal{A}$ : actions.

- 1:  $S_c \leftarrow (\mathcal{O}_b, \mathcal{A}_b)$  ▷ collect  $n_b$  samples by running  $\pi_O$  in environment
- 2:  $\mathcal{M} \leftarrow \emptyset$  ▷ the set  $\mathcal{M}$  collects triples for each tree
- 3: **for**  $i = 0, \dots, I_{max}$  **do**
- 4:   **for**  $j = 1, \dots, T_{max}$  **do**
- 5:      $T_{i,j} \leftarrow \text{train\_tree}(S_c, d)$  ▷ tree  $T_{i,j}$
- 6:      $\mathcal{M}_j \leftarrow (T_{i,j}, \text{eval\_tree}(T_{i,j}, n_{eps}))$  ▷ triple  $\mathcal{M}_j$  ( $\text{eval\_tree}$  returns  $\mathcal{O}_{T_{i,j}}, \mathcal{R}_{T_{i,j}}$ )
- 7:      $\mathcal{M}.\text{append}(\mathcal{M}_j)$
- 8:   **end for**
- 9:    $(T_{i,*}, \mathcal{O}_{T_{i,*}}, \mathcal{R}_{T_{i,*}}) \leftarrow \text{pick\_best}(\mathcal{M})$  ▷ pick triple  $\mathcal{M}_j$  with highest reward  $\mathcal{R}_{T_{i,j}}$
- 10:    $\hat{\mathcal{O}}_{T_{i,*}} \leftarrow \text{choice}(\mathcal{O}_{T_{i,*}}, n_{samp})$  ▷ pick  $n_{samp}$  random observations
- 11:    $S_{new} \leftarrow (\hat{\mathcal{O}}_{T_{i,*}}, \pi_O(\hat{\mathcal{O}}_{T_{i,*}}))$
- 12:    $S_c \leftarrow S_c \cup S_{new}$
- 13: **end for**
- 14: **return**  $T_{I_{max},*}$  ▷ return the best tree from all iterations

---

After each iteration, the generated observations from the so-far best tree are labeled with the oracle's predictions. The tree will likely revisit the region of these observations in the future, and therefore has to know the correct predictions at these points in the observation space. Finally, the selected samples are joined with the previous ones, and the next iteration begins. After a predetermined number of iterations, the overall best tree is returned. A flowchart of the iterative generation of samples for a given tree depth is shown in Figure 1.

A possible variant is to include in line 11 of Algorithm 1 only samples for which the action of the tree  $\mathcal{A}_T$  and of the oracle differ. While this potentially reduces the training set for the DTs, both the oracle's and the DT's predictions are still required for comparison. Since our tests showed no noticeable difference in performance, we do not delve into this variant further.



**Figure 1.** Flowchart of ITER (Algorithm 1)

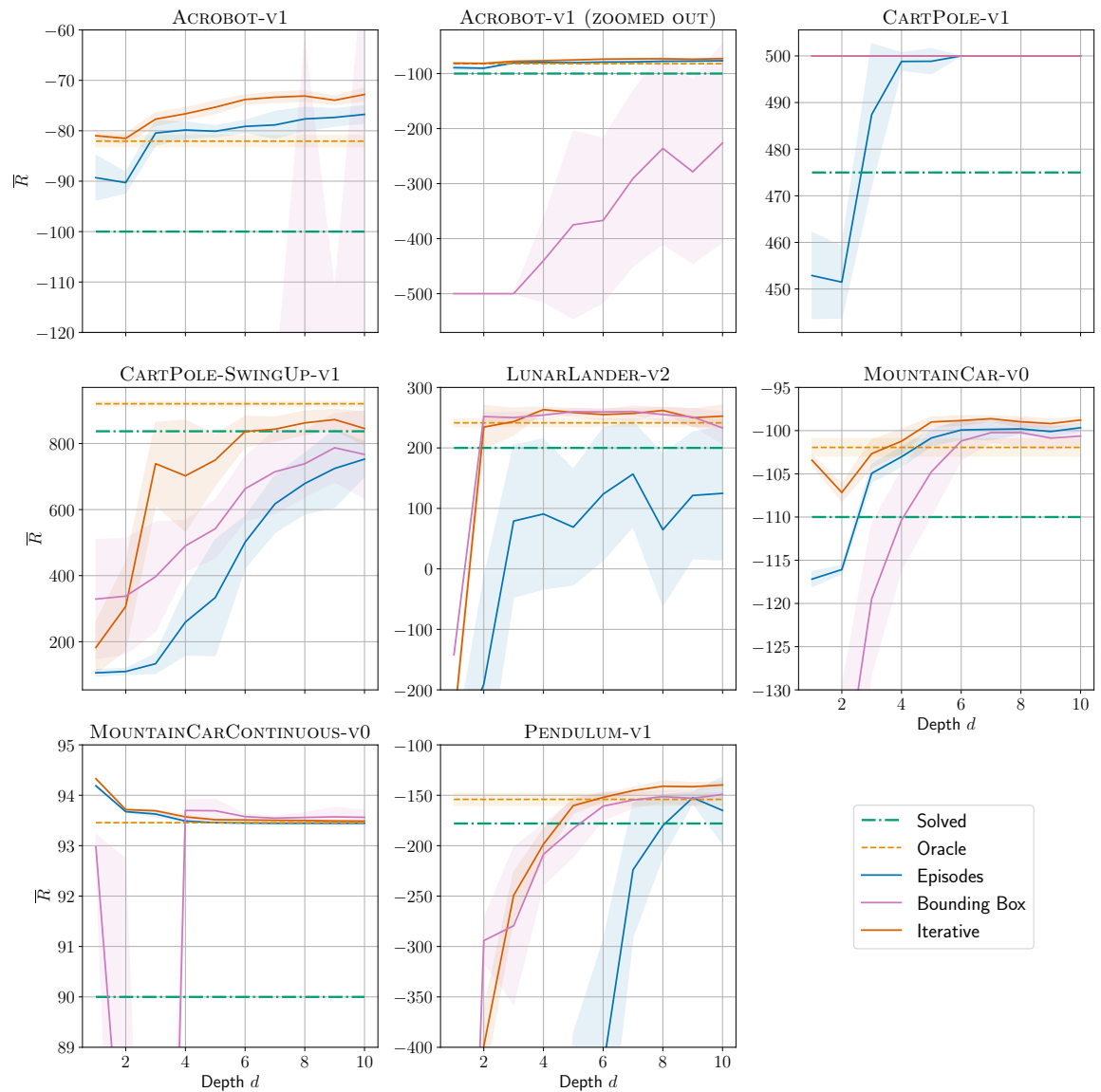
### 3.5. Experimental Setup

All algorithms for sample generation presented in Section 3.4 are evaluated on all environments of Section 3.1 and for all tree depths  $d = 1, \dots, 10$ . Each such trial is called an experiment. For every experiment, the return  $\bar{R}$  is obtained by averaging over 100 evaluation episodes. Each experiment is repeated for 10 runs with different seeds, and both mean  $\mu$  and standard deviation  $\sigma$  of the return  $\bar{R}$  are calculated.

Whenever we create an OPCT, we actually create  $T_{max} = 10$  trees and pick the one with the highest return  $\bar{R}$ . This was done due to the high fluctuation of OPCTs depending on the random seeds that initialized the oblique cuts in the observation space (see implementation in [27]).

In algorithms EPS and BB, we use  $n_s = 30,000$  samples. In ITER (Algorithm 1), we set  $n_b = 20,000$  and  $n_{samp} = 1,000$  to have the same total number of samples across our tested algorithms<sup>2</sup>. The other tunable parameter  $I_{max}$  is set to 10. Adding more than 10 iterations merely resulted in samples generated from a tree that had already been adapted to the oracle. Since these samples hardly differed from those of the oracle, the iterations did not improve the system further.

<sup>2</sup> The ratio between  $n_b$  and  $n_{samp}$  or the total number of samples  $n_s$  are hyperparameters that have not yet been tuned to optimal values. It could be that other values would lead to significantly higher sample efficiency.



**Figure 2.** Return  $\bar{R}$  as a function of DT depth  $d$  for all environments and all presented algorithms. The *solved*-threshold is shown as dash-dotted green line, the average oracle performance as dashed orange line, and the DTs performances as solid lines with average and  $\pm 1\sigma$  of ten repetitions as shaded area. Note: (i) For ACROBOT we show two plots to include BB curve. (ii) Good performance in CARTPOLE leads to overlapping curves: oracle, BB, and ITER are all constantly at  $\bar{R} = 500$ .

## 4. Results

In this section, we discuss results and computational complexity and highlight the advantages of shallow DTs in terms of additional insights and explainability.

### 4.1. Solving Environments

Figure 2 shows our main result: For all environments and all presented algorithms, we plot the return  $\bar{R}$  as a function of DT depth  $d$ . ITER usually reaches the *solved*-threshold (dash-dotted green line) at lower depths than the other two algorithms. It is the only one to solve all investigated environments (at depths up to 10).

Table 2a shows these results in compact quantitative form: Which is the minimum depth needed to reach the *solved*-threshold of an environment?<sup>3</sup> As can be seen, ITER outperforms the other algorithms. Its sum of depths is roughly half of the other algorithms' sums. Depth (1) for MOUNTAINCARCONTINUOUS + BB has been put in brackets because it is somewhat unstable: At depths 2 and 3 the performance is below the threshold, as can be seen from Figure 2. Additionally, Table 2a shows algorithms EPS and ITER for CART as baseline experiments. In both cases, CART does not reach the performance of OPCT, at least not for shallow depths.

Figure 2 shows another remarkable result: DTs can often even surpass the performance of the DRL agents they originate from (the oracles, dashed orange lines). In Table 2b, we report the minimum depths at which DT performances exceeds the oracle's performance. While these numbers can be deceptive too (CARPOLE-SWINGUP + ITER almost reaches oracle performance at depth 8), they show that there is not necessarily a trade-off between model complexity and performance: DTs require orders of magnitude fewer parameters (see Section 4.4.1) and can exhibit higher performance. Again, ITER is the best-performing algorithm in this respect. In Sections 4.3 and 5, we discuss *why* such a performance better than the oracle's can occur.

### 4.2. Computational Complexity

The algorithms presented here differ significantly in computing time, as shown in Table 3. All experiments were performed on a CPU device (i7-10700K  $8 \times 3.8$  GHz). To ensure methodological consistency, we let each algorithm generate the same total number of samples, namely  $n_s = 30,000$ .

In every algorithm, we average over 10 trees for each depth  $d = 1, \dots, 10$  due to fluctuations and represent the time needed in  $t_{OPCT}$ . Additionally, ITER performs  $I_{max} + 1$  iterations for all depths. Hence, a single run takes  $\approx 10 \cdot t_{OPCT}$  time for both EPS and BB while taking  $\approx 10 \cdot (I_{max} + 1) \cdot t_{OPCT}$  time for ITER. (Here, the time  $t_{OPCT}$  reflects the performance of an algorithm in the particular environment, depending on whether more or less time is required for a desired result.) However, the most time-consuming part of a run is the evaluation of the trees, i.e., the generation of samples by running 100 evaluation episodes in the environment. According to these observations, ITER takes the most time to complete.

### 4.3. Decision Space Analysis

Two-dimensional observation spaces offer the possibility of visual inspection of the agents' decision-making. This can provide additional insights and offer explanations of bad results in certain settings. Figure 3 shows the decision surfaces of DRL agents and OPCTs for MOUNTAINCAR in the discrete (top row) and continuous (bottom row) versions. DRL agents can often learn a needlessly complicated partitioning of the observation space, as can be seen prominently in the top left plot, e.g., by yellow patches of "no-acceleration" in the upper left or the blue "accelerate right" area in the bottom right, which are never visited by oracle episodes. This poses a problem for the BB algorithm because it samples, e.g., from

<sup>3</sup> Note that these numbers can be deceptive: A  $10^+$  might be from a tree being below the threshold only by a tiny margin or a 5 might hide the fact that a depth-4 tree is only slightly below the threshold.

**Table 2.** The numbers show the minimum depth of the respective algorithm for **(a) solving an environment** and **(b) surpassing the oracle**. (...) <sup>+</sup> means no DT of investigated depths could reach the threshold. See main text for explanation of **(1)** or **(3)**.

**(a) Solving an environment**

Environment	Model	Algorithms				
		EPS	EPS	BB	ITER	ITER
		CART	OPCT	OPCT	CART	OPCT
Acrobot-v1	DQN	<b>1</b>	<b>1</b>	10 <sup>+</sup>	<b>1</b>	<b>1</b>
CartPole-v1	PPO	3	3	<b>1</b>	3	<b>1</b>
CartPole-SwingUp-v1	DQN	10 <sup>+</sup>	10 <sup>+</sup>	10 <sup>+</sup>	10 <sup>+</sup>	<b>7</b>
LunarLander-v2	PPO	10 <sup>+</sup>	10 <sup>+</sup>	<b>2</b>	10 <sup>+</sup>	<b>2</b>
MountainCar-v0	DQN	3	3	5	3	<b>1</b>
MountainCarContinuous-v0	TD3	<b>1</b>	<b>1</b>	<b>(1)</b>	<b>1</b>	<b>1</b>
Pendulum-v1	TD3	10 <sup>+</sup>	9	6	7	<b>5</b>
Sum		38 <sup>+</sup>	37 <sup>+</sup>	35 <sup>+</sup>	35 <sup>+</sup>	<b>18</b>

**(b) Surpassing the oracle**

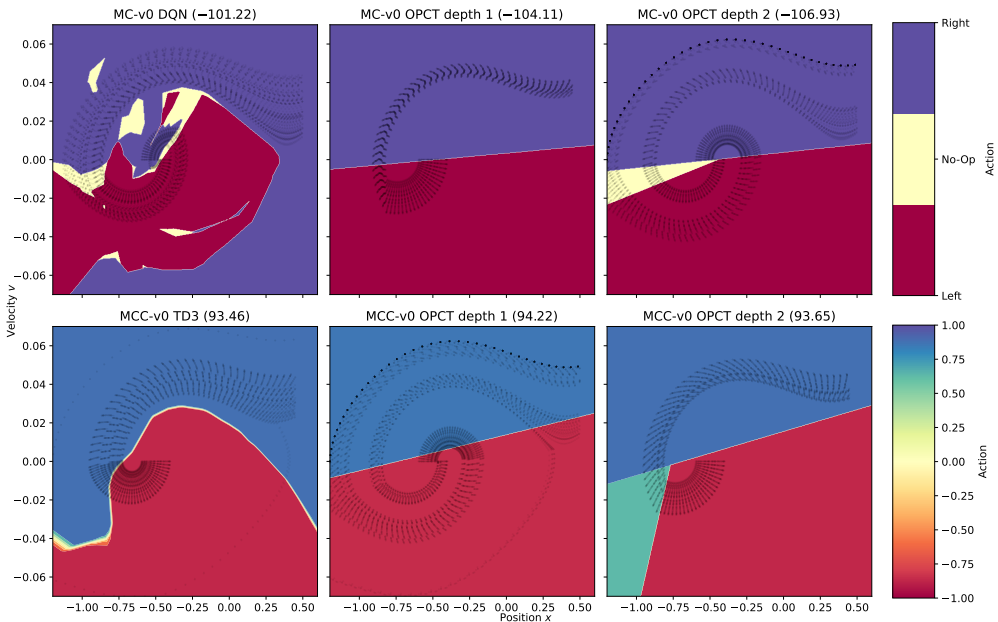
Acrobot-v1	DQN	3	3	10 <sup>+</sup>	2	<b>1</b>
CartPole-v1	PPO	5	6	<b>1</b>	(3)	<b>1</b>
CartPole-SwingUp-v1	DQN	10 <sup>+</sup>	10 <sup>+</sup>	10 <sup>+</sup>	10 <sup>+</sup>	10 <sup>+</sup>
LunarLander-v2	PPO	10 <sup>+</sup>	10 <sup>+</sup>	<b>2</b>	10 <sup>+</sup>	3
MountainCar-v0	DQN	<b>3</b>	5	6	<b>3</b>	4
MountainCarContinuous-v0	TD3	10 <sup>+</sup>	<b>1</b>	4	2	<b>1</b>
Pendulum-v1	TD3	10 <sup>+</sup>	9	8	10	<b>6</b>
Sum		51 <sup>+</sup>	44 <sup>+</sup>	41 <sup>+</sup>	40 <sup>+</sup>	<b>26<sup>+</sup></b>

**Table 3.** Computation times of our algorithms for some environments (we selected those with different observation space dimensions  $D$ ). Shown are the averaged results from 10 runs.  $t_{run}$ : average time elapsed for one run covering all depths  $d = 1, \dots, 10$ .  $t_{OPCT}$ : average time elapsed to train and evaluate 10 OPCTs.  $n_s$ : total number of samples used. Note how the bad performance of BB in the ACROBOT environment leads to longer episodes and therefore longer evaluation times  $t_{OPCT}$ .

Environment	dim $D$	Algorithm	$t_{run}$ [s]	$t_{OPCT}$ [s]	$n_s$
Acrobot-v1	6	EPS	$124.82 \pm 1.96$	$11.71 \pm 0.95$	30,000
		BB	$477.70 \pm 7.89$	$47.56 \pm 2.92$	30,000
		ITER	$1\,189.14 \pm 28.70$	$10.76 \pm 1.02$	30,000
CartPole-v1	4	EPS	$187.39 \pm 3.64$	$18.06 \pm 2.29$	30,000
		BB	$202.65 \pm 1.26$	$19.11 \pm 1.13$	30,000
		ITER	$2\,040.16 \pm 21.46$	$18.50 \pm 1.64$	30,000
MountainCar-v0	2	EPS	$85.39 \pm 1.50$	$7.93 \pm 0.57$	30,000
		BB	$104.92 \pm 2.05$	$10.27 \pm 0.54$	30,000
		ITER	$849.90 \pm 19.79$	$7.68 \pm 0.61$	30,000

the blue bottom right area and therefore, DTs learn in part “wrong” actions, resulting in poor performance. Thus, the inspection of the oracle’s decision space reveals the reason for the partial failure of the BB algorithm in the case of MOUNTAINCAR at low depths.

On the other hand, algorithm ITER will not be trained with samples from regions never visited by oracle or tree episodes. Hence, it has no problem with the blue bottom right area. It delivers more straightforward rules at depth  $d = 1$  or  $d = 2$  (shown in columns 2 and 3 of Figure 3), and yields returns slightly better than the oracle because it generalizes better given its fewer degrees of freedom.



**Figure 3.** Decision surfaces for MOUNTAINCAR (MC, discrete actions) and MOUNTAINCARCONTINUOUS (MCC, continuous actions). Column 1 shows the DRL agents. The OPCTs in column 2 (depth 1) and column 3 (depth 2) were trained with algorithm ITER. The number in brackets in each subplot title is the average return  $\bar{R}$ . The black dots represent trajectories of 100 evaluation episodes with same starting conditions across the different agents. They indicate the regions actually visited in the observation space.

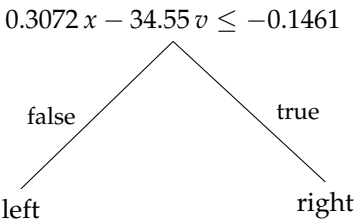
4.4. Explainability for Oblique Decision Trees in RL

Explainability is significantly more difficult to achieve for agents operating in environments with multiple time steps than in those with single time steps (by “single time step” we mean that a decision follows directly in response to a single input record, e.g., classification or regression). This is because the RL agent has to execute a, possibly long, sequence of action decisions before collecting the final return (also known as the credit assignment problem [28,29]).

Given their simple, transparent nature, DTs offer interesting options for explainability in such environments with multiple time steps, which are not immediately applicable to opaque DRL models. This section will illustrate steps towards explainability in RL through shallow DTs.

4.4.1. Decision Trees: Compact and Readable Models

Trees are useful for XAI since they are usually much more compact than DRL models and consist of a set of explicit rules, which are simple in that they are linear inequalities in the input features (for oblique DTs). As an example of the models’ simplicity, Figure 4 shows a very compact OPCT of depth 1 for the MOUNTAINCAR challenge. On the other hand, DRL models have many trainable parameters and form complex, nonlinear features from the input observations. Table 4 shows the number of trainable weights for all our SB3 oracles (DRL models) which are, in all cases, orders of magnitude larger than the number of trainable tree parameters.



**Figure 4.** OPCT of depth 1 obtained by ITER solving MOUNTAINCAR with a return of  $\bar{R} = -107.47$  in 100 evaluation episodes

**Table 4.** Oracle and tree complexity. ( $D$  denotes the dimension of the observation space)

Environment	dim $D$	Model	Number of parameters		depth $d$
			Oracle	ITER	
Acrobot-v1	6	DQN	136,710	9	1
CartPole-v1	4	PPO	9 155	7	1
CartPole-SwingUp-v1	5	DQN	534,534	890	7
LunarLander-v2	8	PPO	9 797	31	2
MountainCar-v0	2	DQN	134,656	5	1
MountainCarContinuous-v0	2	TD3	732,406	5	1
Pendulum-v1	3	TD3	734,806	156	5

The number of parameters in the oracle can be computed from the architecture of the respective policy network. For example, in MOUNTAINCAR with input dimension  $D = 2$ , the DQN consists of 2 networks: a Q- and a target network. Each has  $D + 1 = 2 + 1$  inputs (including bias), 3 action outputs, and a (256, 256) hidden layer architecture<sup>4</sup> which results in

$$2 \cdot ((2 + 1) \cdot 256 + (256 + 1) \cdot 256 + 256 \cdot 3) = 134,656 \text{ weights.}$$

<sup>4</sup> We do not use overly complex DRL architectures but keep the default parameters suggested by the SB3 methods [24].

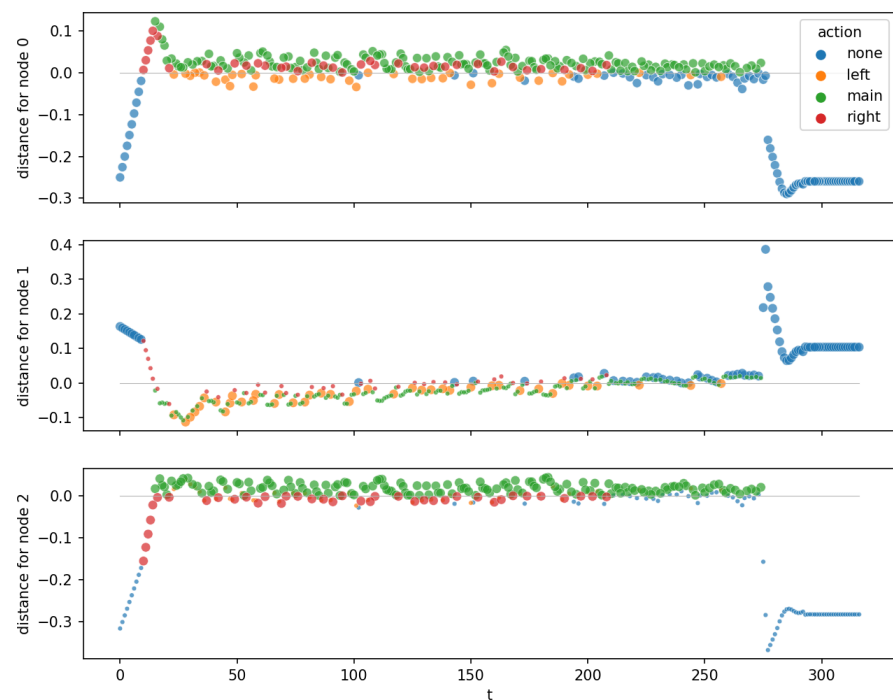
The number of parameters in the tree is a function of the tree's depth and the input dimension: For example, in LUNARLANDER with input dimension  $D = 8$ , the oblique tree of depth  $d = 2$  (as obtained by, e.g., ITER) has  $2^d - 1 = 3 = 1 + 2$  split nodes and  $2^d = 4$  leaf nodes. Each split node has  $D + 1 = 9$  parameters (one weight for each input dimension + threshold), and each leaf node has one adjustable output parameter (the action), resulting in

$$3 \cdot 9 + 4 = 31 \text{ parameters.}$$

However, it should be emphasized that all our DTs have been trained with the guidance of a DRL oracle. At least so far, we could not find a procedure to construct successful DTs solely from interaction with the environment. Optimizing the tree parameters from scratch with respect to cumulative episode reward did not lead to success for tricky OpenAI Gym control problems because incremental changes to a start tree usually miss the goal and thus miss the reward. Only with the guidance of the oracle samples showing reasonable solutions can the tree learn a successful first arrangement of its hyperplanes, which can be refined further.

In a sense, the DT “explains” the oracle by offering a simpler surrogate action model. The surrogate is often nearly as good as or even better than the DRL model in terms of mean reward.

It is a surprising result of our investigations that trees of small depth delivering such high rewards could be found at all. Most prominently, for the environments ACROBOT and LUNARLANDER with their higher input dimensions 6 and 8, it was not expected beforehand to find successful trees of depth 1 and 2, respectively.



**Figure 5.** Distance to hyperplanes for LUNARLANDER. The diagrams show root node, left child, and right child from top to bottom. Small points: node is “inactive” for this observation (the observation is not routed through this node). Big points: node is “active”.

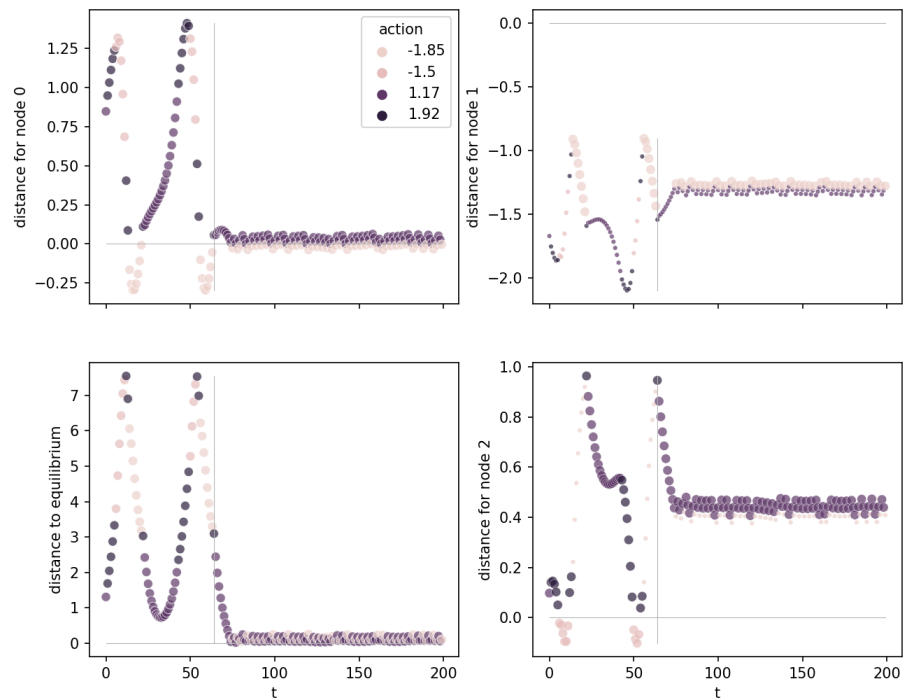
4.4.2. Distance to Hyperplanes

Further insights that help to explain a DT can be gained by visualizing the *distance to hyperplanes* of the nodes of a DT: According to Equation 1, each split node has an associated hyperplane in the observation space so that observations above or below this hyperplane take a route to different subtrees (or leaves) of this node.

Although hyperplanes in higher dimensions are difficult to visualize and interpret, the distance to each hyperplane is easy to visualize, regardless of the dimension of the observation space. It can be shown that even an anisotropic scaling of the observation space leaves the relations between all distances to a particular hyperplane intact, i.e., they are changed only by a common factor.

Figure 5 shows the distance plots for the three nodes 0,1,2 of the ITER OPCT for LUNARLANDER. A distance plot shows the distances of the observations to the node’s hyperplane for a specific episode as a function of time step  $t$ . Each observation is colored by the action the DT has taken. It is seen in Figure 5 that after a short transient period ( $t \leq 40$ ), the distances for every node stay small while the lander sinks. At about time step  $t \approx 280$ , the lander touches the ground and the distances increase (because all velocities are forced to be zero, all positions and angles are forced to certain values prescribed by the shape of the ground), which is, however, irrelevant for the control of the lander and the success of the episode. All nodes are attracting nodes (points are attracted to the respective hyperplane), which bring the ship to an equilibrium (a sinking position with constant  $v_y$ ) that guarantees a safe landing.

Another example is shown in Figure 6 for the three nodes 0,1,2 of the ITER OPCT for PENDULUM. These are the three upper nodes of a tree with depth 3 and, therefore, seven split nodes. Node 0 is an attracting one (“equilibrium node”) because the distances approach zero (here: for  $t \geq 64$ ). Nodes 1 and 2 do not exhibit a zero distance in the equilibrium state; instead, they are “swing-up nodes” responsible for bringing energy into the system through appropriate action switching.



**Figure 6.** Distance plots for PENDULUM. Three of the diagrams show the distance to the hyperplanes of root node 0, left child node 1, and right child node 2. The lower left plot shows the distance to the equilibrium point  $(1, 0, 0)$  in observation space. The vertical line at  $t = 64$  marks the step where the distance for node 0 approaches zero.

The lower left plot (“distance to equilibrium”) of Figure 6 shows the distance of each observation to the unstable equilibrium point  $(\cos \theta, \sin \theta, \omega) = (1, 0, 0)$  (the desired goal state). Interestingly, at  $t = 64$  (vertical gray line), where the distance to node 0 is already close to zero, the distance to the equilibrium point is not yet zero (but soon will be). This is because the hyperplane is oriented in such a way that for small angle  $\theta$ , the angular velocity  $\omega$  is approximately  $\omega = -5.3 \sin(\theta)$  (read off from the node’s parameters). The hyperplane tells us on which path the pole approaches the equilibrium point: The control strategy is to raise  $\omega$  for points below the hyperplane and to lower  $\omega$  for points above. This amounts to the fact that a tilted pole gets an  $\omega$  (corresponding to a certain kinetic energy) such that the pole will come to rest at  $\theta = 0$ . The exact physical equation requires that the kinetic energy has the same value as the necessary change in potential energy,

$$\begin{aligned} E_{kin}(\omega) &= \Delta E_{pot}(\theta) \\ \frac{1}{6} m \ell^2 \omega^2 &= mg \frac{\ell}{2} - mg \frac{\ell}{2} \cos(\theta), \end{aligned}$$

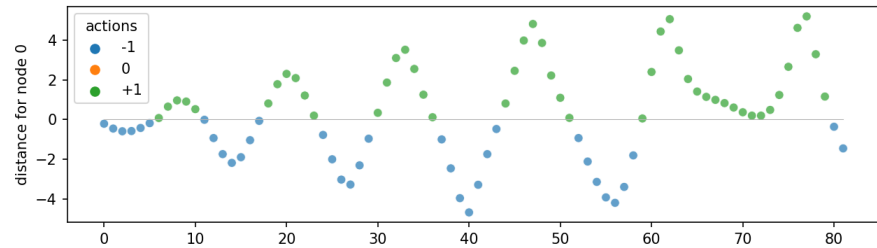
leading to the angular velocity

$$\omega = -\sqrt{\frac{3g}{2\ell}(1 - \cos \theta)} \approx -5.5 \sin(\theta/2) \quad (\text{for } g = 10, \ell = 1)$$

which is not too far from the control strategy found by the node.

Similar patterns (equilibrium nodes + swing-up nodes) can be observed in those environments where the goal is to reach or maintain certain (unstable) equilibrium points (CARTPOLE-SWINGUP, CARTPOLE). Environments like MOUNTAINCAR(CONTINUOUS) or

ACROBOT, which do not have an equilibrium state as goal, consist only of swing-up nodes as shown in Figure 7 for ACROBOT.



**Figure 7.** Distance plot for ACROBOT. The tree has only one node.

In summary, the distance plots show us two characterizing classes of nodes:

- **equilibrium nodes** (or attracting nodes): nodes with distance near to zero for all observation points after a transient period. These nodes occur only in those environments where an equilibrium has to be reached.
- **swing-up nodes**: nodes that do not stabilize at a distance close to zero. These nodes are responsible for swing-up, for bringing energy into a system. The trees for environments MOUNTAINCAR, MOUNTAINCARCONTINUOUS, and ACROBOT consist only of swing-up nodes because, in these environments, the goal state to reach is not an equilibrium.

#### 4.4.3. Sensitivity Analysis

An individual rule of an oblique tree is simple in mathematical terms (it is a linear inequality). It is, nevertheless, difficult to interpret for humans. This is because the individual input features (usually different physical quantities) are multiplied with weights and added up, which has no direct physical interpretation.

For example, the single-node tree of CARPOLE has the rule

$$4.904x - 0.1829v + 19.1497\theta + 1.5203\omega \leq -0.0003 \quad (2)$$

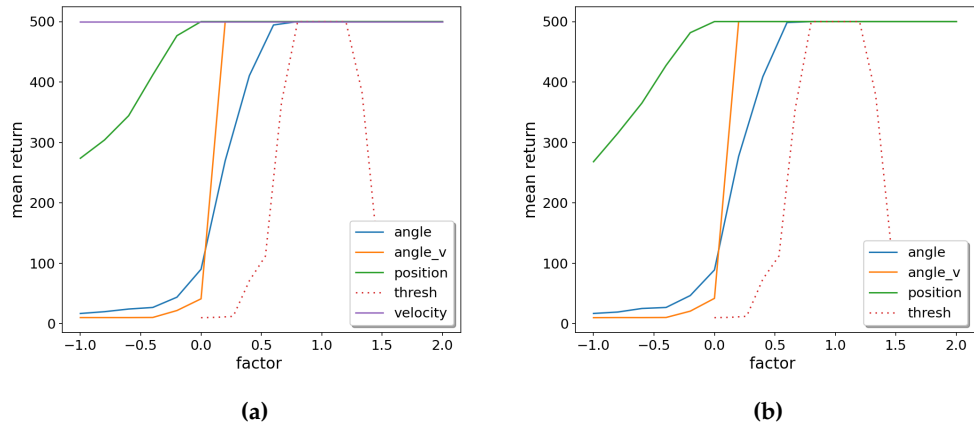
$$\Leftrightarrow w_x x + w_v v + w_\theta \theta + w_\omega \omega \leq \tau$$

which does not directly tell which input features are important for success and is as such difficult to interpret. However, due to the simple mathematical relationship, we may vary each weight  $w_i$  or the threshold  $\tau$  individually and measure the effect on the reward. This is shown in Figure 8. The weights are varied in a wide range of  $[-100\%, 200\%]$  of their nominal value, the threshold  $\tau$  (being close to 0) is varied in the range  $\tau - \bar{w} + [0\%, 200\%] \cdot \bar{w}$ , where  $\bar{w}$  is the mean of all weight magnitudes  $|w_i|$  of the node.

From Figure 8a, we learn that the angle  $\theta$  is the most important input feature: If  $w_\theta$  is below its nominal value, performance starts to degrade, while a higher  $w_\theta$  keeps the optimal reward. The most important parameter is the threshold  $\tau$ , which has to be close to 0; otherwise, performance quickly degrades on both sides. This makes sense because the essential control strategy for the CARPOLE is to react to slightly positive or negative  $\theta$  with the opposite action. On the other hand, the weight for velocity  $v$  is the least important parameter: varying it in Figure 8a does not change the reward, and setting its weight  $w_v$  to 0 and varying the other weights leads to essentially the same picture in Figure 8b.

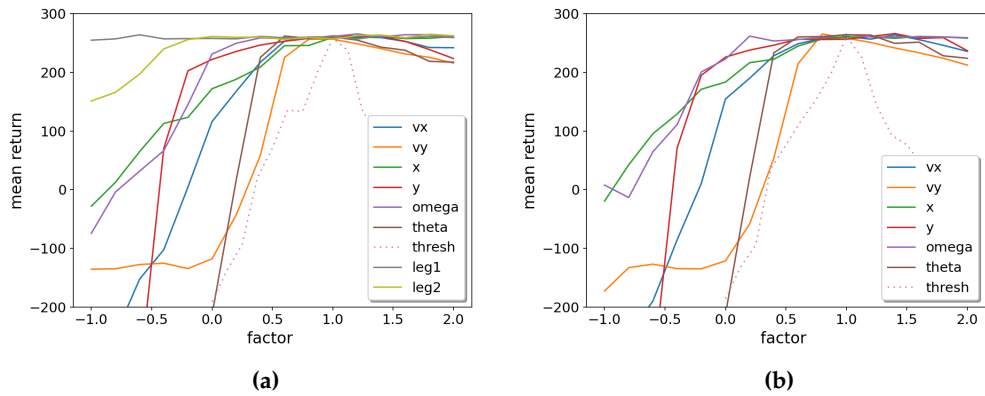
One could object that these sensitivity analysis results could have also been read off directly from Equation 2 because  $\theta$  has the largest, and  $v$  has the most negligible weight in magnitude. However, this is only an accidental coincidence: For many other nodes that we examined using sensitivity analysis, the most/least important feature in terms of mean return did *not* coincide with the largest/smallest weight.

Suppose an OPCT has more than one node, e.g., for LUNARLANDER at depth  $d = 2$  with its three split nodes. In that case, we can perform a sensitivity analysis for each node



**Figure 8.** OPCT sensitivity analysis for CARTPOLE: (a) including all weights, (b) with velocity weight  $w_v$  set to zero. Parameter `angle` is the most important, `velocity` is least important.

separately, or vary the respective weights and thresholds for all nodes simultaneously. The latter is shown in Figure 9. We see that strengthening their weights (setting them to higher values than nominal ones) only leads to slow degradation for all input features. A faster degradation occurs when lowering the weight for a certain input feature. The features  $v_y$ ,  $\theta$ , and  $v_x$  (in this order) are the most important in that respect. On the other hand, the legs (boolean variables signaling the ground-contact of the lander’s legs), are least important: The curves in Figure 9b only change slightly if the legs are removed from the decisions by setting their weights to zero. This makes sense because the legs have a constant *no-contact* value during the majority of episode steps (and when the legs do reach ground-contact, the success or failure of an episode is usually already determined).



**Figure 9.** OPCT sensitivity analysis for LUNARLANDER: (a) including all weights, (b) with the weights for the legs set to zero.  $v_y$  is most important, the ground-contacts of the legs are least important.

Note that the sensitivity analysis shown here is, in a sense, more detailed than a Shapley value analysis [30] because it allows to disentangle the effects of intensifying or weakening an input feature. In contrast, the Shapley value only compares the presence or absence of a feature. (It is of course also less detailed, because it does not take coalitions of features into account as Shapley does.) It is also computationally less intensive: Measuring the Shapley value for RL problems by computing the mean reward over many environment episodes and many coalitions would be, in most cases, prohibitively time-consuming.

5. Discussion

For a broader perspective, we want to point out a number of issues related to our findings.

1. Can we understand *why* ITER + OPCT is successful, while OPCT obtained by the EPS algorithm is not? We have no direct proof, but from the experiments conducted, we can share these insights: Firstly, equilibrium problems like CARPOLE or PENDULUM exhibit a severe oversampling of some states (i.e., the upright pendulum state). If, as a consequence, the initial DT is not well-performing, ITER allows to add the correct oracle actions for problematic samples and to improve the tree in critical regions of the observation space. Secondly, for environment PENDULUM we observed that the initial DT was “bad” because it learned only from near-perfect oracle episodes: It had not seen any samples outside the near-perfect episode region and hypothesized the wrong actions in the *outside* regions. Again, ITER helps: A “bad” DT is likely to visit these *outside* regions and then add samples combined with the correct oracle actions to its sample set.
2. We investigated another possible variant of an iterative algorithm: In this variant, observations were sampled while the oracle was training (and probably also visited unfavorable regions of the observation space). After the oracle had completed its training, the samples were labeled with the action predictions of the fully trained oracle. These labeled samples were presented to DTs in a similar iterative algorithm as described above. However, this algorithm was not successful.
3. A striking feature of our results is that DTs (predominantly those trained with ITER) can outperform the oracle they were trained from. This seems paradoxical at first glance, but the decision space analysis for MOUNTAINCAR in Section 4.3 has shown the likely reason: In Figure 3, the decision spaces of the oracles are overly complex (similar to overfitted classification boundaries, although *overfitted* is not the proper term in the RL context). With their significantly reduced degrees of freedom, the DTs provide simpler, better generalizing models. Since they do not follow the “mistakes” of the overly complex DRL models, they exhibit a better performance. It fits to this picture that the MOUNTAINCARCONTINUOUS results in Figure 2 are best at  $d = 1$  (better than oracle) and slowly degrade towards the oracle performance as  $d$  increases: The DT gets more degrees of freedom and mimics the (slightly) non-optimal oracle better and better. The fact that DTs can be better than the DRL model they were trained on is compatible with the reasoning of Rudin [18], who stated that transparent models are not only easier to explain but often also outperform black box models.
4. We applied algorithm ITER to OPCT. Can CART also benefit from ITER? – The *solved*-depths shown in column ITER + CART of Table 2a add up to  $35^+$ , which is nearly as bad as EPS + CART ( $38^+$ ). Only the *solved*-depth for PENDULUM improved somewhat from  $10^+$  to 7. We conclude that the expressive power of OPCTs is important for being successful in creating shallow DTs with ITER.
5. A last point worth mentioning is the intrinsic trustworthiness of DTs. They partition the observation space in a finite (and often small) set of regions, where the decision is the same for all points. (This feature includes smooth extrapolation to yet-unknown regions.) DRL models, on the other hand, may have arbitrarily complex decision spaces: If such a model predicts action  $a$  for point  $x$ , it is not known which points in the neighborhood of  $x$  will have the same action.

6. Conclusion

In this article we have shown that a considerable range of classic control RL problems can be solved with DTs. Even higher-dimensional problems (ACROBOT and LUNARLANDER) can be solved with surprisingly simple trees. We have presented three algorithms for generating training data for DTs. The training data are composed of the RL environment’s observations and the corresponding DRL oracle’s decisions.

These algorithms differ in the way the points of the observation space are selected: In EPS, all samples visited by the oracle during evaluation episodes are used, BB means that random samples from a hyperrectangle in the observation space are taken, while ITER is an algorithm that takes points in the observation space visited by previous iterations of

trained DTs. Our experiments show how DTs derived with these algorithms can generally not only solve the challenges posed by classic control problems at very moderate depths but also reach or even surpass the performance of the oracle. ITER produces good results on all tested environments with equal or lower DT depths than the other two algorithms. Furthermore, we discuss the advantages of transparent DT models with very few parameters, especially compared to DRL networks, and show how they allow to gain insights that would otherwise be hidden by opaque DRL decision-making processes.

However, our algorithms still require DRL agents as a prerequisite for creating successful DTs.

Future work should test our algorithms on increasingly complex RL environments. ITER still requires fine-tuning of parameters, such as the number of base samples and number of samples per iteration, to optimize sample efficiency and performance.

Our work is intended to provide helpful arguments for simple, transparent RL agents and to advance the knowledge in the field of explainable RL.

**Author Contributions:** Conceptualization, R.C.E., M.L., L.W., and W.K.; investigation, R.C.E., M.O., and W.K.; methodology, R.C.E., M.O., M.L., L.W., and W.K., project administration, R.C.E., M.O., and W.K.; software, R.C.E., M.O., and W.K.; supervision, L.W. and W.K.; validation, R.C.E., M.O., and W.K.; visualization, R.C.E., M.O., and W.K.; writing—original draft, R.C.E., M.O., and W.K.; writing—review & editing, R.C.E., M.O., M.L., L.W., and W.K.;

**Funding:** This research was supported by the research training group “Dataninja” (Trustworthy AI for Seamless Problem Solving; Next Generation Intelligence Joins Robust Data Analysis) funded by the German federal state of North Rhine-Westphalia.

**Institutional Review Board Statement:** Not applicable

**Data Availability Statement:** The experiments presented in this article are openly available in our Github repository <https://github.com/MarcUnknown/Iterative-Oblique-Decision-Trees>

**Conflicts of Interest:** The authors declare no conflict of interest. The funders had no role in the design of the study; in the collection, analyses, or interpretation of data; in the writing of the manuscript; or in the decision to publish the results.

Abbreviations

The following abbreviations are used in this manuscript:

BB	Bounding Box algorithm
CART	Classification and Regression Trees
DRL	Deep Reinforcement Learning
DQN	Deep Q-Network
DT	Decision Tree
EPS	Episode Samples algorithm
ITER	Iterative Training of Explainable RL models
OPCT	Oblique Predictive Clustering Tree
PPO	Proximal Policy Optimization
RL	Reinforcement Learning
SB3	Stable-Baselines3
TD3	Twin Delayed Deep Deterministic Policy Gradient

References

1.

Engelhardt, R.C.; Lange, M.; Wiskott, L.; Konen, W. Sample-Based Rule Extraction for Explainable Reinforcement Learning. In Proceedings of the Machine Learning, Optimization, and Data Science; Nicosia, G.; Ojha, V.; La Malfa, E.; La Malfa, G.; Pardalos, P.; et al., Eds.; Springer: Cham, 2023; Vol. 13810, LNCS, pp. 330–345. [https://doi.org/10.1007/978-3-031-25599-1\\_25](https://doi.org/10.1007/978-3-031-25599-1_25).

473

2.

Adadi, A.; Berrada, M. Peeking inside the black-box: a survey on explainable artificial intelligence (XAI). *IEEE access* **2018**, *6*, 52138–52160. <https://doi.org/10.1109/ACCESS.2018.2870052>.

476

3.

Molnar, C.; Casalicchio, G.; Bischl, B. Interpretable Machine Learning – A Brief History, State-of-the-Art and Challenges. In Proceedings of the ECML PKDD 2020 Workshops; Koprinska, I.; Kamp, M.; Appice, A.; et al., Eds.; Springer: Cham, 2020; pp. 417–431. [https://doi.org/10.1007/978-3-030-65965-3\\_28](https://doi.org/10.1007/978-3-030-65965-3_28).

478

4. Molnar, C. Interpretable Machine Learning: A guide for making black box models explainable. <https://christophm.github.io/interpretable-ml-book/>, 2022. 481
5. Puiutta, E.; Veith, E.M.S.P. Explainable Reinforcement Learning: A Survey. In Proceedings of the Machine Learning and Knowledge Extraction; Holzinger, A.; Kieseberg, P.; Tjoa, A.M.; Weippl, E., Eds.; Springer: Cham, 2020; pp. 77–95. [https://doi.org/10.1007/978-3-030-57321-8\\_5](https://doi.org/10.1007/978-3-030-57321-8_5). 482
6. Heuillet, A.; Couthouis, F.; Díaz-Rodríguez, N. Explainability in deep reinforcement learning. *Knowledge-Based Systems* **2021**, 214, 106685. <https://doi.org/10.1016/j.knosys.2020.106685>. 483
7. Milani, S.; Topin, N.; Veloso, M.; Fang, F. A survey of explainable reinforcement learning, 2022. <https://doi.org/10.48550/arXiv.2202.08434>. 484
8. Lundberg, S.M.; Erion, G.; Chen, H.; DeGrave, A.; Prutkin, J.M.; Nair, B.; Katz, R.; Himmelfarb, J.; Bansal, N.; Lee, S.I. From local explanations to global understanding with explainable AI for trees. *Nature Machine Intelligence* **2020**, 2, 56–67. <https://doi.org/10.1038/s42256-019-0138-9>. 485
9. Liu, G.; et al. Toward Interpretable Deep Reinforcement Learning with Linear Model U-Trees. In Proceedings of the Machine Learning and Knowledge Discovery in Databases; Berlingerio, M.; et al., Eds. Springer, 2019, Vol. 11052, LNCS, pp. 414–429. [https://doi.org/10.1007/978-3-030-10928-8\\_25](https://doi.org/10.1007/978-3-030-10928-8_25). 486
10. Mania, H.; Guy, A.; Recht, B. Simple random search of static linear policies is competitive for reinforcement learning. In Proceedings of the Advances in Neural Information Processing Systems; Bengio, S.; Wallach, H.; Larochelle, H.; Grauman, K.; Cesa-Bianchi, N.; Garnett, R., Eds. Curran Associates, Inc., 2018, Vol. 31. 487
11. Coppens, Y.; Efthymiadis, K.; et al. Distilling deep reinforcement learning policies in soft decision trees. In Proceedings of the IJCAI 2019 workshop on explainable artificial intelligence, 2019, pp. 1–6. 488
12. Frosst, N.; Hinton, G.E. Distilling a Neural Network Into a Soft Decision Tree. In Proceedings of the First International Workshop on Comprehensibility and Explanation in AI and ML; Besold, T.R.; Kutz, O., Eds., 2017, Vol. 2071, CEUR Workshop Proceedings. 489
13. Verma, A.; Murali, V.; Singh, R.; Kohli, P.; Chaudhuri, S. Programmatically Interpretable Reinforcement Learning. In Proceedings of the 35th International Conference on Machine Learning; Dy, J.; Krause, A., Eds. PMLR, 2018, Vol. 80, *Proceedings of Machine Learning Research*, pp. 5045–5054. 490
14. Zilke, J.R.; Loza Mencía, E.; Janssen, F. DeepRED – Rule Extraction from Deep Neural Networks. In Proceedings of the Discovery Science; Calders, T.; Ceci, M.; Malerba, D., Eds.; Springer: Cham, 2016; Vol. 9956, LNCS, pp. 457–473. [https://doi.org/10.1007/978-3-319-46307-0\\_29](https://doi.org/10.1007/978-3-319-46307-0_29). 491
15. Qiu, W.; Zhu, H. Programmatic Reinforcement Learning without Oracles. In Proceedings of the International Conference on Learning Representations, 2022. 492
16. Schapire, R.E. The strength of weak learnability. *Machine Learning* **1990**, 5, 197–227. <https://doi.org/10.1007/BF00116037>. 493
17. Freund, Y.; Schapire, R.E. A Decision-Theoretic Generalization of On-Line Learning and an Application to Boosting. *Journal of Computer and System Sciences* **1997**, 55, 119–139. <https://doi.org/10.1006/jcss.1997.1504>. 494
18. Rudin, C. Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead. *Nature Machine Intelligence* **2019**, 1, 206–215. <https://doi.org/10.1038/s42256-019-0048-x>. 495
19. Brockman, G.; Cheung, V.; Pettersson, L.; Schneider, J.; Schulman, J.; Tang, J.; Zaremba, W. OpenAI Gym, 2016. <https://doi.org/10.48550/arXiv.1606.01540>. 496
20. Lovatto, A.G. CartPole Swingup - A simple, continuous-control environment for OpenAI Gym. <https://github.com/0xangelo/gym-cartpole-swingup>, 2021. 497
21. Schulman, J.; Wolski, F.; Dhariwal, P.; Radford, A.; Klimov, O. Proximal Policy Optimization Algorithms, 2017. <https://doi.org/10.48550/arXiv.1707.06347>. 498
22. Mnih, V.; Kavukcuoglu, K.; et al. Playing Atari with Deep Reinforcement Learning, 2013. <https://doi.org/10.48550/arXiv.1312.5602>. 499
23. Fujimoto, S.; van Hoof, H.; Meger, D. Addressing Function Approximation Error in Actor-Critic Methods. In Proceedings of the 35th International Conference on Machine Learning; Dy, J.; Krause, A., Eds. PMLR, 2018, Vol. 80, *Proceedings of Machine Learning Research*, pp. 1587–1596. 500
24. Raffin, A.; Hill, A.; Gleave, A.; Kanervisto, A.; Ernestus, M.; Dormann, N. Stable-Baselines3: Reliable Reinforcement Learning Implementations. *Journal of Machine Learning Research* **2021**, 22, 1–8. 501
25. Breiman, L.; Friedman, J.H.; Olshen, R.A.; Stone, C.J. *Classification And Regression Trees*; Routledge, 1984. 502
26. Pedregosa, F.; Varoquaux, G.; Gramfort, A.; Michel, V.; Thirion, B.; Grisel, O.; Blondel, M.; Prettenhofer, P.; Weiss, R.; Dubourg, V.; et al. Scikit-learn: Machine Learning in Python. *Journal of Machine Learning Research* **2011**, 12, 2825–2830. 503
27. Stepišnik, T.; Kocev, D. Oblique predictive clustering trees. *Knowledge-Based Systems* **2021**, 227, 107228. <https://doi.org/10.1016/j.knosys.2021.107228>. 504
28. Alipov, V.; Simmons-Edler, R.; Putintsev, N.; Kalinin, P.; Vetrov, D. Towards practical credit assignment for deep reinforcement learning, 2021. <https://doi.org/10.48550/arXiv.2106.04499>. 505
29. Woergoetter, F.; Porr, B. Reinforcement learning. *Scholarpedia* **2008**, 3, 1448. <https://doi.org/10.4249/scholarpedia.1448>. 506
30. Roth, A.E., Ed. *The Shapley Value: Essays in Honor of Lloyd S. Shapley*; Cambridge University Press, 1988. <https://doi.org/10.1017/CBO9780511528446>. 507

**Disclaimer/Publisher’s Note:** The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.

539

540

541