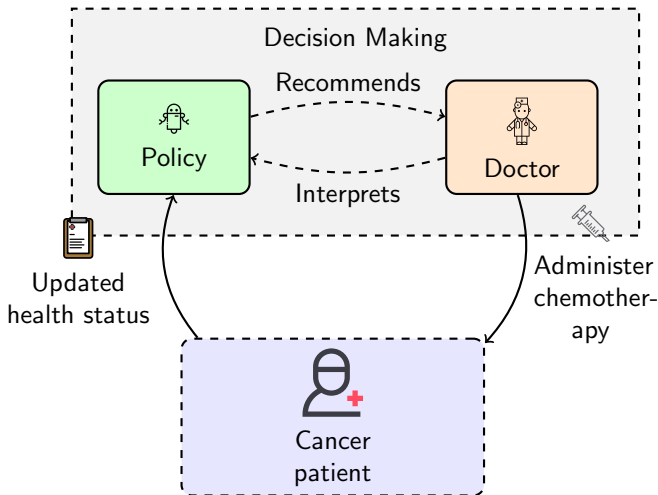# Interpretability, Decision Trees, and Sequential Decision Making

Hector Kohler

Supervised by Dr. Riad Akrour (HdR) and Prof. Philippe Preux (HdR)
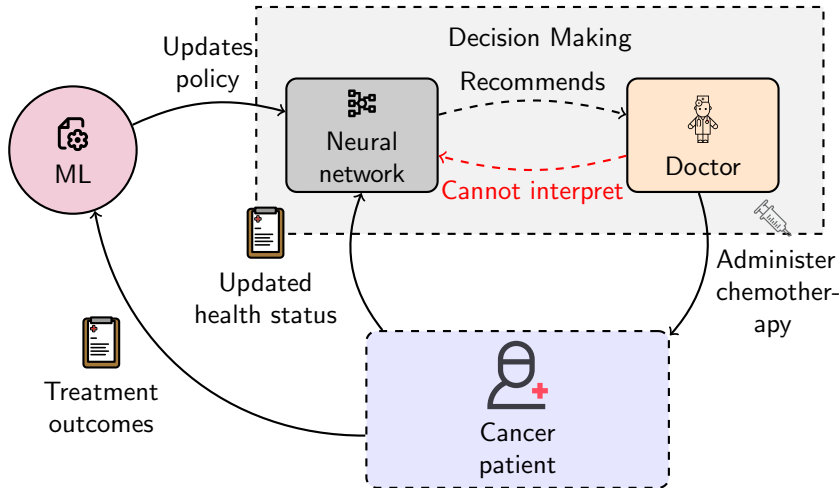Université de Lille, CNRS, Inria, UMR CRIStAL 9189, France

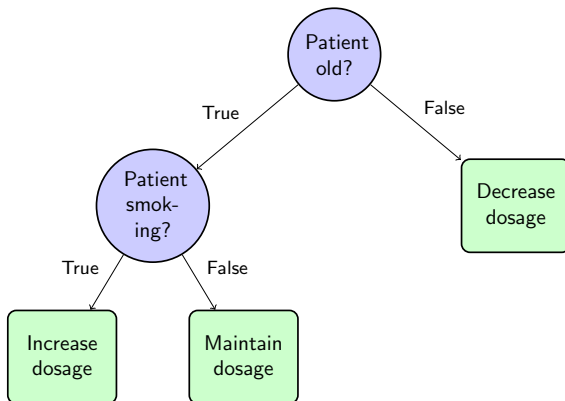November 22, 2025

# Sequential decision making (SDM)



Sequential decision making in cancer treatment.

# Machine learning (ML) of policies for SDM



Machine learning of neural networks has many recent successes but neural networks are black-box.
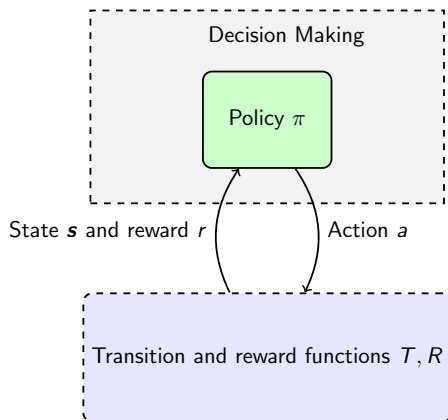
# Decision trees



A generic decision tree of depth $D = 2$.

Successful algorithms for non-sequential learning: (Bertsimas and Dunn 2017; Breiman et al. 1984; Demirovic et al. 2022; Mazumder, Meng, and Wang 2022; Verwer and Zhang 2019) . . . What about SDM?

# Markov decision processes (MDPs) and reinforcement learning (RL)



Decision Making

Policy $\pi$

State $s$ and reward $r$    Action $a$

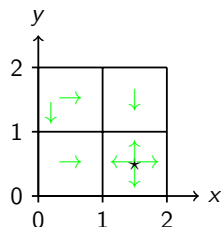Transition and reward functions $T, R$

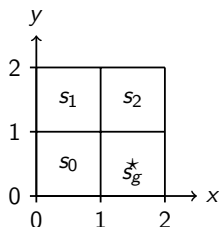Markov decision processes (Puterman 1994).

- RL (Sutton and Barto 1998) aims to find a policy, $\pi : S \rightarrow A$ that maximizes:

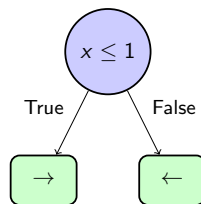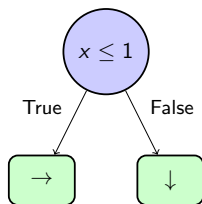$$\mathop{\mathbb{E}}_{s_t \sim T} \left[ \sum_{t=0}^{\infty} \gamma^t R(s_t, \pi(s_t)) \right]$$

- Lots of succesful RL algorithms (Mnih et al. 2015; Schulman et al. 2017; Sutton and Barto 1998).

- No interpretability concerns.

# Grid world MDP



A grid world MDP and optimal actions.



An optimal depth-1 decision tree policy and a sub-optimal depth-1 decision tree policy.

# Indirect approach: imitation learning



**Step 1:** Use NN to generate states

**Step 2:** Use NN to obtain actions

**Step 3:** Use supervised learning to train a decision tree

Imitation learning works well in practice to get interpretable policies (Bastani, Pu, and Solar-Lezama 2018; Milani et al. 2024; Ross, Gordon, and Bagnell 2010) but no optimality guarantees.

# Contributions

1. Why is learning optimal interpretable policies for sequential decision making difficult?

2. How to leverage sequential decision making to learn interpretable classifiers for supervised learning?

3. How to measure policy interpretability in sequential decision making?

# Iterative bounding Markov decision processes (IBMDP)



Trajectory in an IBMDP of the grid world MDP (Topin et al. 2021). Actions build a decision tree policy and rewards control the interpretability-performance trade-off.

# Pros and cons of IBMDPs

## Pros

- No need to design new algorithm: we can use deep RL.
- IBMDP rewards trade-off naturally interpretability and performances.

## Cons

- Only **determinstic** and **partially observable** (a.k.a. memoryless or reactive) policies are equivalent to decision tree policies.
- Finding the best **deterministic** and **partially observable** policy is NP-hard (Littman 1994)!

## Re-formulation

*Q: Can we use reinforcement learning to directly optimize trade-offs of performance and interpretability in SDM?*

$$\Longleftrightarrow$$

*Q: How does RL perform for optimizing **deterministic** and **partially observable** policies in IBMDPs?*

# Result: RL cannot retrieve optimal depth-1 trees for the grid world MDP



Distributions of final tree policies learned with various (asymmetric) RL algorithms (Baisero and Amato 2022; Baisero, Daley, and Amato 2022; Loch and Singh 1998; Singh, Jaakkola, and Jordan 1994; Sutton and Barto 1998) across 100 seeds. For each different performance-interpretability trade-off value $\zeta$, each point represent the share of different trees.

# Result: for similar problems, RL struggles when there is partial observability (not surprising)



Success rates of different (asymmetric) RL algorithms over thousands of runs when applied to learning either deterministic partially observable policies in an IBMDP deterministic Markovian policies in the same IBMDP.

Classification MDP and the unique optimal depth-1 tree.

**We show that deterministic partially observable policies for classification IBMDPs ($\Leftrightarrow$ decision tree policies) are in fact Markovian.**

# Result: RL can retrieve optimal depth-1 trees for the toy classification MDPs



Distributions of final tree policies learned with various RL algorithms across 100 seeds. For each different performance-interpretability trade-off value $\zeta$, each point represent the share of different trees.

# Perspectives for direct RL of decision tree policies.

- Interpretability for SDM problems can be difficult because of **partial observability**.
- Should we focus on indirect approaches? Hybrid approaches (Wu et al. 2020)?
- Fixing the policy tree structure a priori (paramteric trees, Marton et al. 2025)?
- Design algorithms that learn deterministic partially observable policies (Lambrechts, Bolland, and Ernst 2025; Lambrechts, Ernst, and Mahajan 2025)?

## RL works in classification MDPs

*Q: Can we leverage SDM design new decision tree induction algorithms for the supervised learning setting?* **A: Yes!**

# Decision trees in supervised learning

- $N$ data points. Each $x_i$ is described by $p$ features and has a label $y_i \in \mathcal{Y}$.

$$\mathcal{L}(T) = \frac{1}{N}\sum_{i=1}^{N}\ell(y_i, f(x_i)) + \alpha C(T)$$

- Trees **interpretable** and **competitive with neural nets** (Grinsztajn, Oyallon, and Varoquaux 2022).

- Greedy algorithms **sub-optimal accuracy**, but $O(2^D)$ operations (Breiman et al. 1984; J Ross Quinlan 1993; J. R. Quinlan 1986) .

- Optimal algorithms, **optimal accuracy**, but $O((2Np)^D)$ operations (NP-hard) (Bertsimas and Dunn 2017; Chaouki, Read, and Bifet 2024; Demirovic et al. 2022; Hyafil and Rivest 1976; Linden, Weerdt, and Demirović 2023).

Pathological Dataset (N=10_000, p=2)

CART - Acc: 75.88% Operations: 7, Time: 0.008s

Optimal - Acc: 100% Operations: 173, Time: 4s

DPDT (Ours) - Acc: 100.0% Operations: 9, Time: 0.04s

A checkers board data set highlights the limitations of existing works.

# Decision tree induction as solving MDPs

### Intuition

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

- S: data subsets.
- A: test or leaf nodes that can be added to the tree.
- R: penalty or accuracies.
- T: node traversals.

MDP formulation of a generic decision tree induction for a supervised learning task.

- Greedy algorithms consider only one candidate action in each state which is the test that minimizes some impurity criterion → **MDP state space size is $O(2^D)$**.

- Optimal algorithms consider all possible actions in each state → **MDP state space size is $O((2Np)^D)$**.

- Let's choose candidate actions adaptively → for each MDP state consider $B$ actions: **state space size is $O((2B)^D)$**.

Overview of our algorithm DPDT presented at the 31st ACM SIGKDD conference.

---

[1]Because states are entire datasets, we implement DPDT with a depth-first search to limit the space complexity.

# Comparing tree accuracy to complexity

Train accuracy and operation count when learning depth-3 decision trees.

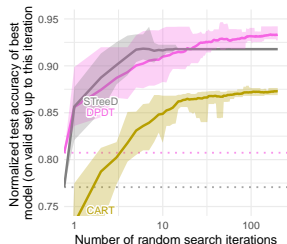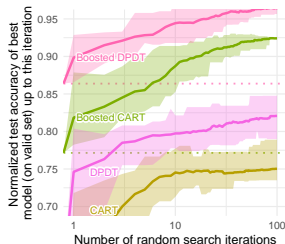| | | | Accuracy | | | | | Operations | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Dataset | N | p | Opt Quant-BnB | Greedy CART | DPDT light | DPDT full | | Opt Quant-BnB | Greedy CART | DPDT light | DPDT full |
| room | 8103 | 16 | **0.992** | 0.968 | 0.991 | **0.992** | | $10^6$ | 15 | 286 | 16100 |
| bean | 10888 | 16 | **0.871** | 0.777 | 0.812 | 0.853 | | $5 \cdot 10^6$ | 15 | 295 | 25900 |
| eeg | 11984 | 14 | **0.708** | 0.666 | 0.689 | 0.706 | | $2 \cdot 10^6$ | 13 | 289 | 26000 |
| avila | 10430 | 10 | **0.585** | 0.532 | 0.574 | **0.585** | | $3 \cdot 10^7$ | 9 | 268 | 24700 |
| magic | 15216 | 10 | **0.831** | 0.801 | 0.822 | 0.828 | | $6 \cdot 10^6$ | 15 | 298 | 28000 |
| htru | 14318 | 8 | **0.981** | 0.979 | 0.979 | 0.980 | | $6 \cdot 10^7$ | 15 | 295 | 25300 |
| occup. | 8143 | 5 | **0.994** | 0.989 | 0.991 | **0.994** | | $7 \cdot 10^5$ | 13 | 280 | 16300 |
| skin | 196045 | 3 | **0.969** | 0.966 | 0.966 | 0.966 | | $7 \cdot 10^4$ | 15 | 301 | 23300 |
| fault | 1552 | 27 | **0.682** | 0.553 | 0.672 | 0.674 | | $9 \cdot 10^8$ | 13 | 295 | 24200 |
| segment | 1848 | 18 | **0.887** | 0.574 | 0.812 | 0.879 | | $2 \cdot 10^6$ | 7 | 220 | 16300 |
| page | 4378 | 10 | **0.971** | 0.964 | 0.970 | 0.970 | | $10^7$ | 15 | 298 | 22400 |
| bidding | 5056 | 9 | **0.993** | 0.981 | 0.985 | **0.993** | | $3 \cdot 10^5$ | 13 | 256 | 9360 |
| raisin | 720 | 7 | **0.894** | 0.869 | 0.879 | 0.886 | | $4 \cdot 10^6$ | 15 | 295 | 20900 |
| rice | 3048 | 7 | **0.938** | 0.933 | 0.934 | 0.937 | | $2 \cdot 10^7$ | 15 | 298 | 25500 |
| wilt | 4339 | 5 | **0.996** | 0.993 | 0.994 | 0.995 | | $3 \cdot 10^5$ | 13 | 274 | 11300 |
| bank | 1097 | 4 | **0.983** | 0.933 | 0.971 | 0.980 | | $6 \cdot 10^4$ | 13 | 271 | 7990 |

DPDT depth-5 trees vs. other detph-5 trees

Boosted DPDT vs. Boosted CART

Boosted DPDT vs. other classifiers

## Perspectives

- New SOTA decision tree induction with dynamic programming in MDPs.
- What about using DPDT for indirect decision tree policy learning for SDM?
- What performances could we reach with an industry-grade implementation of XGboost+DPDT?

### Let us take a step back

*Q: Are decision trees really the most interpretable model?*
**A: It depends.**

# Policy interpretability

Policy performance



High

Neural networks

Tree ensembles

Linear policies

Decision trees

Low

Low                                    High    Policy interpretability

**Heuristic** interpretability-performance trade-offs of different policy classes.
Interpretability is often presented in opposition to performances.

# How to measure policy interpretability?

## Challenges (Doshi-Velez and Kim 2017; Glanois et al. 2024; Lipton 2018)

- No definition of interpretability.
- Measuring might require humans.

## The notion of *simulatability* (Lipton 2018)

- Interpretability $\simeq$ how long for human to make the same computations.
- Interpretability $\simeq$ how much effort for a human to read through the entire policy.
- Less parameters mean more interpretability (Freitas 2014; Lavrač 1999).
- Time to formally verify a policy decreases with interpretability (Barceló et al. 2020).

# A methodology to measure policy interpretability without humans

## Simulatability (Lipton 2018)

1. How long it takes for human to make the same computations given an input $\simeq$ policy inference time.
2. How much effort it would take a human to read through the entire policy once $\simeq$ policy size in memory.

## Not that simple in practice (Luo et al. 2024)

- Different hardwares (CPUs vs GPUs).
- Different implementations (matrix operations vs fully sequentially) . . .

# We propose policy unfolding

```
# Decision tree for Mountain Car
def play(x):
    if x[1] <= -0.2597:
        if x[1] <= -0.6378:
            return 0
        else:
            if x[0] <= -1.0021:
                return 2
            else:
                return 0
    else:
        if x[1] <= -0.0508:
            if x[0] <= 0.2979:
                if x[0] <= 0.0453:
                    return 2
                else:
                    if x[1] <=
    -0.2156:
                        return 0
                    else:
                        return 2
            else:
                return 0
        else:
            return 2
```
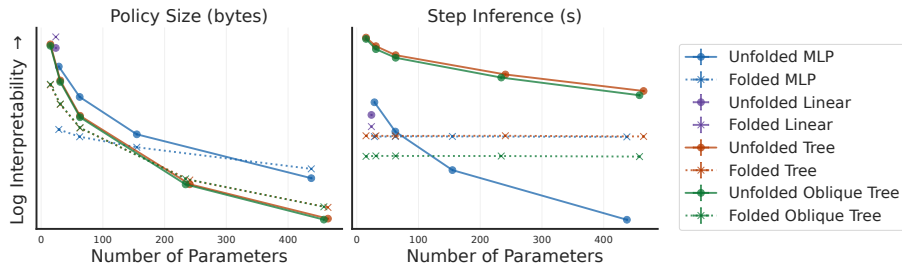
```
# Small ReLU MLP for Pendulum
def play(x):
    h_layer_0_0 = 1.238*x[0]+0.971*x
        [1]
                      +0.430*x[2]+0.933
    h_layer_0_0 = max(0, h_layer_0_0
        )
    h_layer_0_1 = -1.221*x[0]+1.001
                      *x[1]-0.423*x[2]
                      +0.475
    h_layer_0_1 = max(0, h_layer_0_1
        )
    h_layer_1_0 = -0.109*h_layer_0_0
                      -0.377*h_layer_0_1
                      +1.694
    h_layer_1_0 = max(0, h_layer_1_0
        )
    h_layer_1_1 = -3.024*h_layer_0_0
                      -1.421*h_layer_0_1
                      +1.530
    h_layer_1_1 = max(0, h_layer_1_1
        )

    h_layer_2_0 = -1.790*h_layer_1_0
                      +2.840*h_layer_1_1
                      +0.658
    y_0 = h_layer_2_0
    return [y_0]
```

1. Does our methodology respect consensus on policy interpretability?
2. Is policy unfolding necessary?
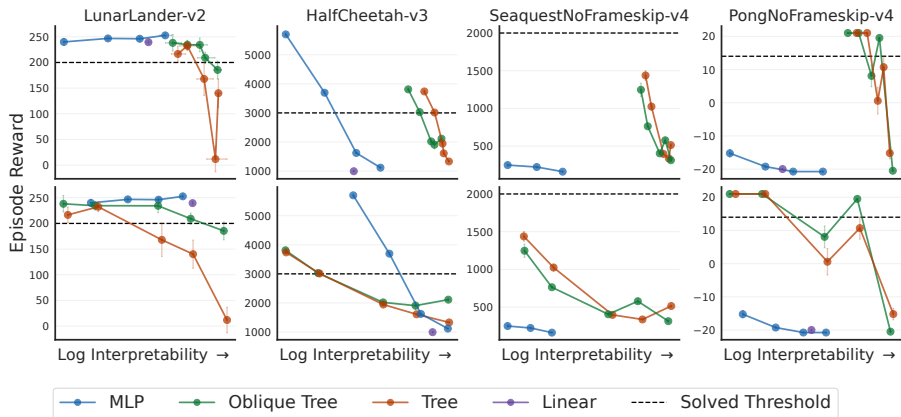3. What kind of results we can obtain using our proposed methodology?

### W

e imitate $\sim 40000$ expert policies from `stable-baselines3` using various policy classes/nb parameters on various environments.

Aggregated policies interpretability on classic control environments

# Result: there is no dominating policy class for all environments



Interpretability-Performance trade-offs. Top row, interpretability is measured with step inference times. Bottom row, the interpretability is measured with policy size.

# Perspectives

- Beliefs such as "trees are more interpretable than neural networks" should be used with caution.
- Tree-like policy classes can have good inductive bias (e.g. Atari).
- Can a human study confirm our results?
- What about (very) big models?
- Can we use our policy programs as low level skills (hierarchical RL)?

# Conclusion: interpretable machine learning is a difficult research topic

- Technical challenges: **partial observability in SDM, NP-hardness**.
  $\rightarrow$ Focus on indirect approaches and/or on POMDP research first.
- Fundamental challenges: **no definition**.
  $\rightarrow$ Discuss with the community (InterpPol workshop).
- **Decision trees offer good inductive bias for SDM in games or tabular data**.

### My hope

Motivate interpretability by finding a real-world problem where interpretability is *really* necessary (Nagendran et al. 2024).

📄 Baisero, Andrea and Christopher Amato (2022). "Unbiased Asymmetric Reinforcement Learning under Partial Observability". In: *Proceedings of the 21st International Conference on Autonomous Agents and Multiagent Systems*. AAMAS '22. Virtual Event, New Zealand: International Foundation for Autonomous Agents and Multiagent Systems, pp. 44–52. ISBN: 9781450392136.

📄 Baisero, Andrea, Brett Daley, and Christopher Amato (Jan. 2022). "Asymmetric DQN for partially observable reinforcement learning". In: *Proceedings of the Thirty-Eighth Conference on Uncertainty in Artificial Intelligence*. Ed. by James Cussens and Kun Zhang. Vol. 180. Proceedings of Machine Learning Research. PMLR, pp. 107–117. URL: https://proceedings.mlr.press/v180/baisero22a.html.

📄 Barceló, Pablo et al. (2020). "Model interpretability through the lens of computational complexity". In: *Advances in neural information processing systems*.

📄 Bastani, Osbert, Yewen Pu, and Armando Solar-Lezama (2018). "Verifiable Reinforcement Learning via Policy Extraction". In.

📄 Bertsimas, Dimitris and Jack Dunn (2017). "Optimal classification trees". In: *Machine Learning* 106, pp. 1039–1082.

📄 Breiman, L et al. (1984). *Classification and Regression Trees*. Wadsworth.

📄 Chaouki, Ayman, Jesse Read, and Albert Bifet (2024). "Branches: A Fast Dynamic Programming and Branch & Bound algorithm for Optimal Decision Trees". In: arXiv: 2406.02175 [cs.LG]. URL: https://arxiv.org/abs/2406.02175.

📄 Demirovic, Emir et al. (2022). "MurTree: Optimal Decision Trees via Dynamic Programming and Search". In: *Journal of Machine Learning Research* 23.26, pp. 1–47. URL: http://jmlr.org/papers/v23/20-520.html.

📄 Doshi-Velez, Finale and Been Kim (2017). "Towards A Rigorous Science of Interpretable Machine Learning". In: arXiv: 1702.08608 [stat.ML]. URL: https://arxiv.org/abs/1702.08608.

📄 Freitas, Alex A. (Mar. 2014). "Comprehensible classification models: a position paper". In: *SIGKDD Explor. Newsl.* 15.1, pp. 1–10. ISSN: 1931-0145. DOI: 10.1145/2594473.2594475. URL: https://doi.org/10.1145/2594473.2594475.

📄 Glanois, Claire et al. (2024). "A survey on interpretable reinforcement learning". In: *Machine Learning*, pp. 1–44.

📄 Grinsztajn, Léo, Edouard Oyallon, and Gaël Varoquaux (2022). "Why do tree-based models still outperform deep learning on typical tabular data?" In: *Advances in neural information processing systems* 35, pp. 507–520.

📄 Hyafil, Laurent and Ronald L. Rivest (1976). "Constructing optimal binary decision trees is NP-complete". In: *Information Processing Letters* 5.1, pp. 15–17. ISSN: 0020-0190. DOI: https://doi.org/10.1016/0020-0190(76)90095-8. URL: https://www.sciencedirect.com/science/article/pii/0020019076900958.

📄 Lambrechts, Gaspard, Adrien Bolland, and Damien Ernst (2025). "Informed POMDP: Leveraging Additional Information in Model-Based RL". In: *Reinforcement Learning Journal* 2, pp. 763–784.

📄 Lambrechts, Gaspard, Damien Ernst, and Aditya Mahajan (2025). "A Theoretical Justification for Asymmetric Actor-Critic algorithms". In: *Forty-second International Conference on Machine Learning*. URL: https://openreview.net/forum?id=F1yANMCnAn.

📄 Lavrač, Nada (1999). "Selected techniques for data mining in medicine". In: *Artificial Intelligence in Medicine* 16.1. Data Mining

Techniques and Applications in Medicine, pp. 3–23. ISSN: 0933-3657.
DOI: https://doi.org/10.1016/S0933-3657(98)00062-1. URL:
https://www.sciencedirect.com/science/article/pii/
S0933365798000621.

📄 Linden, Jacobus van der, Mathijs de Weerdt, and Emir Demirović
(2023). "Necessary and Sufficient Conditions for Optimal Decision
Trees using Dynamic Programming". In: *Advances in Neural
Information Processing Systems* 36. Ed. by A. Oh et al.,
pp. 9173–9212.

📄 Lipton, Zachary C. (2018). "The Mythos of Model Interpretability: In
machine learning, the concept of interpretability is both important and
slippery.". In: *Queue* 16.3, pp. 31–57.

📄 Littman, Michael L. (1994). "Memoryless policies: theoretical
limitations and practical results". In: *Proceedings of the Third
International Conference on Simulation of Adaptive Behavior: From
Animals to Animats 3: From Animals to Animats 3*. SAB94. Brighton,
United Kingdom: MIT Press, pp. 238–245. ISBN: 0262531224.

📄 Loch, John and Satinder P. Singh (1998). "Using Eligibility Traces to
Find the Best Memoryless Policy in Partially Observable Markov

Decision Processes". In: *Proceedings of the Fifteenth International Conference on Machine Learning*. ICML '98. San Francisco, CA, USA: Morgan Kaufmann Publishers Inc., pp. 323–331. ISBN: 1558605568.

📄 Luo, Lirui et al. (2024). "End-to-End Neuro-Symbolic Reinforcement Learning with Textual Explanations". In: *International Conference on Machine Learning (ICML)*.

📄 Marton, Sascha et al. (2025). "Mitigating Information Loss in Tree-Based Reinforcement Learning via Direct Optimization". In: URL: https://openreview.net/forum?id=qpXctF2aLZ.

📄 Mazumder, Rahul, Xiang Meng, and Haoyue Wang (17–23 Jul 2022). "Quant-BnB: A Scalable Branch-and-Bound Method for Optimal Decision Trees with Continuous Features". In: *Proceedings of the 39th International Conference on Machine Learning*. Proceedings of Machine Learning Research 162. Ed. by Kamalika Chaudhuri et al., pp. 15255–15277. URL: https://proceedings.mlr.press/v162/mazumder22a.html.

📄 Milani, Stephanie et al. (Apr. 2024). "Explainable Reinforcement Learning: A Survey and Comparative Review". In: *ACM Comput. Surv.*

56.7. ISSN: 0360-0300. DOI: 10.1145/3616864. URL: https://doi.org/10.1145/3616864.

📄 Mnih, Volodymyr et al. (2015). "Human-level control through deep reinforcement learning". In: *nature* 518.7540, pp. 529–533.

📄 Nagendran, Myura et al. (2024). "Eye tracking insights into physician behaviour with safe and unsafe explainable AI recommendations". In: *NPJ Digital Medicine* 7.1, p. 202.

📄 Puterman, Martin L. (1994). *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. John Wiley & Sons.

📄 Quinlan, J Ross (1993). "C4. 5: Programs for machine learning". In: *Morgan Kaufmann google schola* 2, pp. 203–228.

📄 Quinlan, J. R. (1986). "Induction of Decision Trees". In: *Mach. Learn.* 1.1, pp. 81–106.

📄 Ross, Stéphane, Geoffrey J. Gordon, and J. Andrew Bagnell (2010). "A Reduction of Imitation Learning and Structured Prediction to No-Regret Online Learning". In.

📄 Schulman, John et al. (2017). "Proximal policy optimization algorithms". In: *arXiv preprint arXiv:1707.06347*.

📄 Singh, Satinder P., Tommi S. Jaakkola, and Michael I. Jordan (1994). "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. New Brunswick, NJ, USA: Morgan Kaufmann Publishers Inc., pp. 284–292. ISBN: 1558603352.

📄 Sutton, Richard S. and Andrew G. Barto (1998). *Reinforcement Learning: An Introduction*. Cambridge, MA: The MIT Press.

📄 Topin, Nicholay et al. (2021). "Iterative bounding mdps: Learning interpretable policies via non-interpretable methods". In: *Proceedings of the AAAI Conference on Artificial Intelligence* 35, pp. 9923–9931.

📄 Verwer, Sicco and Yingqian Zhang (2019). "Learning optimal classification trees using a binary linear program formulation". In: *Proceedings of the AAAI conference on artificial intelligence* 33, pp. 1625–1632.

📄 Wu, Mike et al. (Apr. 2020). "Regional Tree Regularization for Interpretability in Deep Neural Networks". In: 34, pp. 6413–6421. DOI: 10.1609/aaai.v34i04.6112. URL: https://ojs.aaai.org/index.php/AAAI/article/view/6112.