

**UNIVERSITÉ DE LILLE**  
**INRIA**

École doctorale École Gradué MADIS-631

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

Thèse présentée par **Hector KOHLER**

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

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

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

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

Thèse dirigée par Philippe PREUX directeur  
Riad AKROUR co-directeur

**Composition du jury**

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

## COLOPHON

Mémoire de thèse intitulé « Interprétabilité via l'Apprentissage Supervisé ou par Renforcement d'Arbres de Décisions », écrit par Hector KOHLER, achevé le 5 juin 2025, composé au moyen du système de préparation de document  $\text{\LaTeX}$  et de la classe yathesis dédiée aux thèses préparées en France.

UNIVERSITÉ DE LILLE  
INRIA

École doctorale École Gradué MADIS-631

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

Thèse présentée par **Hector KOHLER**

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

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

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

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

Thèse dirigée par Philippe PREUX directeur  
Riad AKROUR co-directeur

**Composition du jury**

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



**UNIVERSITÉ DE LILLE**  
**INRIA**

Doctoral School **École Gradué MADIS-631**

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

Thesis defended by **Hector KOHLER**

Defended on **December 1, 2025**

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

Academic Field **Computer Science**

Speciality **Computer Science and Applications**

# Interpretability through Supervised or Reinforcement Learning of Decision Trees

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

## Committee members

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



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

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Ut purus elit, vestibulum ut, placerat ac, adipiscing vitae, felis. Curabitur dictum gravida mauris. Nam arcu libero, nonummy eget, consectetur id, vulputate a, magna. Donec vehicula augue eu neque. Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Mauris ut leo. Cras viverra metus rhoncus sem. Nulla et lectus vestibulum urna fringilla ultrices. Phasellus eu tellus sit amet tortor gravida placerat. Integer sapien est, iaculis in, pretium quis, viverra ac, nunc. Praesent eget sem vel leo ultrices bibendum. Aenean faucibus. Morbi dolor nulla, malesuada eu, pulvinar at, mollis ac, nulla. Curabitur auctor semper nulla. Donec varius orci eget risus. Duis nibh mi, congue eu, accumsan eleifend, sagittis quis, diam. Duis eget orci sit amet orci dignissim rutrum. Nam dui ligula, fringilla a, euismod sodales, sollicitudin vel, wisi. Morbi auctor lorem non justo. Nam lacus libero, pretium at, lobortis vitae, ultricies et, tellus. Donec aliquet, tortor sed accumsan bibendum, erat ligula aliquet magna, vitae ornare odio metus a mi. Morbi ac orci et nisl hendrerit mollis. Suspendisse ut massa. Cras nec ante. Pellentesque a nulla. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus. Aliquam tincidunt urna. Nulla ullamcorper vestibulum turpis. Pellentesque cursus luctus mauris.

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

---

---

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

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

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

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

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

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

---



# Sommaire

Résumé	vii
Sommaire	ix
Preliminary Concepts	1
<b>I A difficult problem : Learning Decision Trees for MDP</b>	<b>11</b>
1 A Decision Tree Policy for an MDP is a Policy for some Partially Observable MDP	13
2 An attempt at Learning Decision Tree Policies with Reinforcement Learning	15
3 Conclusion	25
<b>II An easier problem : Learning Decision Trees for MDPs that are Classification tasks</b>	<b>27</b>
4 DPDT-intro	29
5 DPDT-paper	31
6 Conclusion	33
<b>III Beyond Decision Trees : what can be done with other Interpretable Policies?</b>	<b>35</b>
Conclusion générale	37

<b>A Programmes informatiques</b>	<b>39</b>
<b>Table des matières</b>	<b>41</b>

# Preliminary Concepts

## Interpretable Sequential Decision Making

### What is Sequential Decision Making?

In this manuscript we are mostly interested in sequential decision making. Humans engage in sequential decision making in all aspects of life. In medicine, doctors have to decide when to use chemotherapy next based on the patient's current health in order to heal (cite). In agriculture, agronomists have to decide when to fertilize next based on the current soil and weather conditions in order to maximize plant growth (cite). In automotive, the auto-pilot system has to decide how to steer the wheel next based on lidar sensors in order to maintain a safe trajectory (cite). In video games, a bot decides what attack to throw next based on the player's and its own state in order to provide the best entertainment (cite). Those sequential decision making processes exhibits key similarities : an agent takes actions based on some current information to achieve some goal. As computer scientists, we ought to design computer programs (cite) that can help humans during those sequential decision making processes. For example, a doctor could benefit from a program that would recommend the "best" treatment given the patient's state. Machine learning algorithms (cite) output such helpful programs. Every day, hundreds of new machine learning algorithms are published<sup>1</sup>. While most of those scientific articles focus on finding the "best" program possible, not many work make sure that their recommendations can be understood by humans. Next, we describe the notion of interpretability that

---

1. <https://arxiv.org/list/cs.LG/pastweek?skip=0&show=2000>

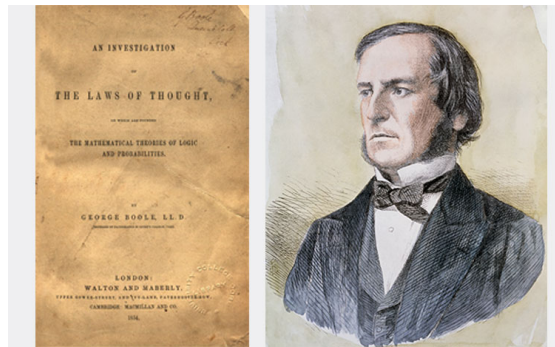


FIGURE 1 – British logician and philosopher George Boole (1815-1864) next to its book *The Laws of Thoughts* (1854) that is the oldest known record of the word “interpretability”.

is key to ensure safe deployment of computer programs trained with machine learning in critical sectors like medicine.

## What is Interpretability?

(figure) Interpretability is a crucial topic in modern science. While computer programs trained with machine learning algorithms have become more and more performing and made their way into our society, they are often black-box : their outputs, e.g., the animals on the image (cite), tokamak control (cite), or even the abstract of your next article (cite), are computed with operations that are too complex for humans to fully understand. Indeed most of recent machine learning breakthroughs are obtained by training very large programs like neural networks (cite) which are black-box by definition since they perform sometimes millions of operations on their inputs. Originally, the etymology of “interpretability” is the Latin “interpretabilis” meaning “that can be understood and explained”.

According to the Oxford English dictionary, the first recorded use of the english word “interpretability” dates back to 1854 when the british logician George Boole (figure) described the addition of concepts :

I would remark in the first place that the generality of a method in Logic must very much depend upon the generality of its elementary

processes and laws. We have, for instance, in the previous sections of this work investigated, among other things, the laws of that logical process of addition which is symbolized by the sign  $+$ . Now those laws have been determined from the study of instances, in all of which it has been a necessary condition, that the classes or things added together in thought should be mutually exclusive. The expression  $x + y$  seems indeed uninterpretable, unless it be assumed that the things represented by  $x$  and the things represented by  $y$  are entirely separate; that they embrace no individuals in common. And conditions analogous to this have been involved in those acts of conception from the study of which the laws of the other symbolical operations have been ascertained. The question then arises, whether it is necessary to restrict the application of these symbolical laws and processes by the same conditions of interpretability under which the knowledge of them was obtained. If such restriction is necessary, it is manifest that no such thing as a general method in Logic is possible. On the other hand, if such restriction is unnecessary, in what light are we to contemplate processes which appear to be uninterpretable in that sphere of thought which they are designed to aid? [(cite)]

What is remarkable is that the supposedly first recorded occurrence of “interpretability” was in the context of (pre-)computer science. Boole asked : *when can we meaningfully apply formal mathematical operations beyond the specific conditions under which we understand them?* In Boole’s era, the concern was whether logical operations like addition could be applied outside their original interpretable contexts—where symbols and their sum represent concepts that humans can understand, e.g.  $\text{red} + \text{apples} = \text{red apples}$ . Today, we face an analogous dilemma with machine learning algorithms : neural networks learn complex un-intelligible combinations of inputs (representations), but we often deploy them in contexts where operations should be understood by humans, e.g., in medicine.

Circling back to our cancer treatment example, we would ideally want doctors to have access to computer programs that can recommend “good” treatments and which operations can be understood. Those two aspects of machine learning

TABLEAU 1 – Related work in *global* interpretable machine learning following the direct/indirect classification. Each class is further divided into the type of interpretable model they output. Methods specific to supervised learning are colored in **red** and methods specific to reinforcement learning are in **blue**

Direct				indirect			
(cite)	(cite)	(cite)	(cite)	(cite)	(cite)	(cite)	(cite)

models–performance and interpretability–often compromise; highly performing models like neural networks are often less interpretable and vice-versa (cite). However, we will observe later on that aiming for interpretability is not necessarily always constraining but can be a quite positive bias for performances in some domains like video games. Interestingly, one of the key challenges of doing research in interpretability is the lack of formalism; there is no definition of what is an interpretable compute program such as an MDP policy. Throughout this manuscript we make the hypothesis that interpretability is the (space and time) complexity of a program (cite) and hence mostly focus on decision trees (low complexity) (cite) and neural networks (high complexity). Despite this lack of formalism the necessity of deploying interpretable models has sparked many works that we present next.

## What are existing approaches for learning interpretable programs?

(figure)(table) Techniques associated with interpretability in machine learning are either global or local. Global methods output a whole model that is interpretable while local approaches LIME FI (cite) output explanations of parts of the model such as how the model considers each part of a single specific input. Inside global interpretable machine learning, approaches are either direct or indirect. Direct algorithms like decision tree induction (cite) are algorithms that directly search a space of interpretable models. On the other hand are indirect methods–sometimes called post-hoc–that attempt to interpret some parts of non-interpretable model (cite). An ideal solution to interpretability would be to have global methods that outputs highl performing models. However there is

no definite answer as to which is better between direct and indirect algorithms. In Table 1 we present both supervised and reinforcement learning works on global interpretability. For reviews of other aspects of interpretable machine learning such as understanding not only the model but also the learning process or on local approaches see (cite). In addition to classifying interpretable machine learning methods by their algorithms type, we can also look at the type of model outputted (tree vs neural network vs linear models...)

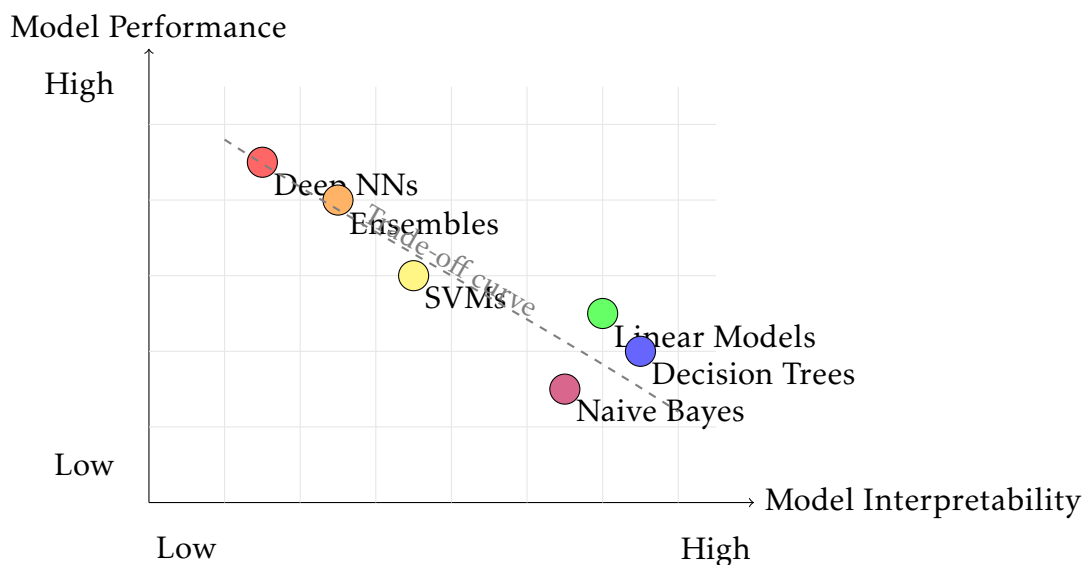


FIGURE 2 – The interpretability-performance trade-off in machine learning. Different model classes are positioned according to their typical interpretability and performance characteristics. The dashed line illustrates the general trade-off between these two properties.

In Figure 2 we present the popular trade-off between interpretability and performance of different model classes.

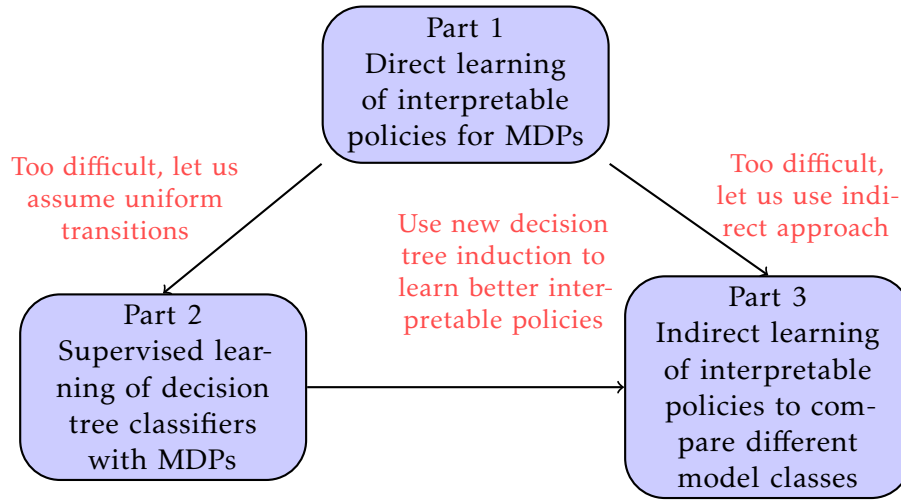


FIGURE 3 – Thesis structure showing the progression from direct reinforcement learning of decision tree policies (Chapter 1) to simplified approaches : supervised learning with uniform transitions (Chapter 2) and indirect learning methods (Chapter 3).

## Outline of the Thesis

### Technical Preliminaries

#### What are decision trees?

(figure) As mentionned earlier, as opposed to neural networks, decision trees are supposidly very interpretable because the only apply boolean operations on the program input without relying on internal complex representations.

**Definition 1** (Decision tree). *A decision tree is a rooted tree  $T = (V, E)$  where :*

- *Each internal node  $v \in V$  is associated with a test function  $f_v : \mathcal{X} \rightarrow \{0, 1\}$  that maps input features  $x \in \mathcal{X}$  to a boolean.*
- *Each edge  $e \in E$  from an internal node corresponds to an outcome of the associated test function.*
- *Each leaf node  $\ell \in V$  is associated with a prediction  $y_\ell \in \mathcal{Y}$ , where  $\mathcal{Y}$  is the output space.*
- *For any input  $x \in \mathcal{X}$ , the tree defines a unique path from root to leaf, determi-*





FIGURE 4 – The american statistician Leo Breiman (1928-2005) author of *Classification and Regression Trees* (1984)

ning the prediction  $T(x) = y_\ell$  where  $\ell$  is the reached leaf.

For non-sequential decision processes such as classification ; it is possible to efficiently compute decision trees with induction over some training examples (cite). However

## How to learn decision trees?

What about reinforcement learning?

## Markov decision processes/problems

(figure) Markov decision processes (MDPs) were first introduced in the 1950s by Richard Bellman (cite). Informally, an MDP models how an agent acts over time to achieve its goal. At every timestep, the agent observes its current state, e.g. a patient weight and tumor size, and takes an action, e.g. injects a certain amount of chemotherapy. When doing a certain action in a certain state, the agent gets a reward that helps it evaluate the quality of its action with respect to its goal, e.g., the tumor size decrease when the agent has to cure cancer. Finally, the agent is provided with a new state, e.g. the updated patient state, and repeats this process over time. Following Martin L. Puterman's book on MDPs (cite), we formally define as follows.

**Definition 2** (Markov decision process). *An MDP is a tuple  $\mathcal{M} = \langle S, A, R, T, T_0 \rangle$  where :*

- *$S$  is a finite set of states  $s \in \mathbb{R}^n$  representing all possible configurations of the environment.*
- *$A$  is a finite set of actions  $a \in \mathbb{Z}^m$  available to the agent.*
- *$R : S \times A \rightarrow \mathbb{R}$  is the reward function that assigns a real-valued reward to each state-action pair.*
- *$T : S \times A \rightarrow \Delta(S)$  is the transition function that maps state-action pairs to probability distributions over next states, where  $\Delta(S)$  denotes the probability simplex over  $S$ .*
- *$T_0 \in \Delta(S)$  is the initial distribution over states.*

Now we can also model the “goal” of the agent. Informally, the goal of an agent is to behave such that it gets as much reward as it can over time. For example, in the cancer treatment case, the best reward the agent can get is to completely get rid of the patient’s tumor after some time. Furthermore, we want our agent to prefer behaviour that gets rid of the patient’s tumor as fast as possible. We can formally model the agent’s goal as an optimization problem as follows.

**Definition 3** (Markov decision problem). *Given an MDP  $\mathcal{M} = \langle S, A, R, T, T_0 \rangle$ , the goal of an agent following policy  $\pi : S \rightarrow A$  is to maximize the expected discounted sum of rewards :*

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

where  $\gamma \in (0, 1)$  is the discount factor that controls the trade-off between immediate and future rewards.

Hence, algorithms presented in this manuscript aim to find solutions to Markov decision problems, i.e. the optimal policy :  $\pi^* = \underset{\pi}{\operatorname{argmax}} J(\pi)$  For the rest of this text, we will use an abuse of notation and denote both a Markov decision process and the associated Markov decision problem by MDP.

## Exact solutions for Markov decision problems

It is possible to compute the exact optimal policy  $\pi^*$  using dynamic programming (cite). Indeed, one can leverage the Markov property to find for all states the best action to take based on the reward of upcoming states.

**Definition 4** (Value of a state). *The value of a state  $s \in S$  under policy  $\pi$  is the expected discounted sum of rewards starting from state  $s$  and following policy  $\pi$  :*

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

Applying the Markov property gives a recursive definition of the value of  $s$  under policy  $\pi$  :

$$V^\pi(s) = R(s, \pi(s)) + \gamma \sum_{s' \in S} T^{s'}(s, \pi(s)) V^\pi(s')$$

where  $T^{s'}(s, \pi(s))$  is the probability of transitioning to state  $s'$  when taking action  $\pi(s)$  in state  $s$ .

**Definition 5** (Optimal value of a state). *The optimal value of a state  $s \in S$ ,  $V^*(s)$ , is the value of state  $s$  when following the optimal policy :  $V^{\pi^*}(s)$ .*

$$V^*(s) = V^{\pi^*}(s) = \max_{\pi} [J(\pi)]$$

Hence, the algorithms we study in the thesis can also be seen as solving the problem :  $\pi^* = \underset{\pi}{\operatorname{argmax}} \mathbb{E}[V^\pi(s_0) \mid s_0 \sim T_0]$ . The well-known Value Iteration algorithm (algorithm) solves this problem exactly (cite).

More realistically, neither the transition kernel  $T$  nor the reward function  $R$  of the MDP are known, e.g., the doctor can't **know** how the tumor and the patient health will change after a dose of chemotherapy, it can only **observe** the change. This distinction between the information available to the agent is paralleled with the distinction between dynamic programming and reinforcement learning (RL) that we describe next.



FIGURE 5 – The godfathers of sequential decision making. Andrew Barto and Richard Sutton are the ACM Turing Prize 2024 laureate and share an advisor advisee relationship.

## Reinforcement learning of approximate solutions to MDPs

(figure) Reinforcement learning algorithms popularized by Richard Sutton (figure) (cite) don't **compute** an optimal policy but rather **learn** an approximate one based on sequences of observations  $(s_t, a_t, r_t, s_{t+1})_t$ . RL algorithms usually fall into two categories : value-based (cite) and policy gradient (cite). The first group of RL algorithms computes an approximation of  $V^*$  using temporal different learning (algorithms) (cite). The other class of RL algorithms leverage the policy gradient theorem (algorithm) (cite) to approximate  $\pi^*$ . Both class of algorithms are known to converge to the optimal value or policy under some conditions (cite) and have known great successes in real-world applications (cite). The books from Puterman, Bertsekas, Sutton and Barto, offer a great overview of MDPs and algorithm to solve them. There are many other ways to learn policies such as simple random search (cite) or model-based reinforcement learning. However, not many algorithms consider the learning of policies that can be easily understood by humans which we discuss next and that is the core of this manuscript.

## **Première partie**

**A difficult problem : Learning  
Decision Trees for MDP**



I have not failed. I've just found  
10.000 ways that won't work.

---

Thomas A. Edison

# Chapitre 1

A Decision Tree Policy for an MDP is a  
Policy for some Partially Observable  
MDP

## **1.1 How to Learn a Decision Tree Policy for an MDP?**

### **1.1.1 Imitation**

### **1.1.2 Soft Trees**

### **1.1.3 Iterative Bounding MDPs**

## **1.2 How to solve Iterative Bounding MDPs?**

### **1.2.1 Asymmetric Reinforcement Learning**

### **1.2.2 Learning a decision tree policy is solving a POMDP**

## **1.3 Is it hard to properly learn a Decision Tree Policy for an MDP?**

### **1.3.1 POMDPs are way harder to solve than MDPs**

### **1.3.2 Memoryless approaches to solve POMDPs seem ineffective**



# An attempt at Learning Decision Tree Policies with Reinforcement Learning

## 2.1 Grid Worlds

$$\begin{aligned}
 V(g) &= \zeta \sum_{i=0}^{\infty} \gamma^{2i} + \sum_{i=0}^{\infty} \gamma^{2i+1} \\
 V(0) &= \zeta + \gamma^2 V(g) \\
 V(1) &= \zeta + \gamma^2 V(0) \\
 \frac{1}{4} \gamma \frac{1}{1-\gamma} + \frac{1}{4} \frac{1}{1-\gamma} &\leq \frac{1}{4} V(g) + \frac{2}{4} V(0) + \frac{1}{4} V(1) \\
 \zeta \cdot \sum_{i=0}^{\infty} \gamma^i &\leq \frac{1}{4} V(g) + \frac{2}{4} V(0) + \frac{1}{4} V(1) \\
 \frac{1}{4} V(g) + \frac{1}{4} (\zeta + \gamma V(0)) + \frac{1}{4} (\zeta + \gamma V(1)) + \frac{1}{4} V(0) &\leq \frac{1}{4} V(g) + \frac{2}{4} V(0) + \frac{1}{4} V(1) \\
 \frac{1}{4} V(g) + \frac{1}{4} V(0) + \frac{1}{4} (\zeta + \gamma^2 \zeta \sum_{i=0}^{\infty} \gamma^{2i}) + \frac{1}{4} (\zeta \sum_{i=0}^{\infty} \gamma^{2i}) &\leq \frac{1}{4} V(g) + \frac{2}{4} V(0) + \frac{1}{4} V(1)
 \end{aligned}$$

### 2.1.1 Step-by-step derivation of the lower bound on $\zeta$

**Step 1 : Simplify the left side of the inequality**

$$\frac{1}{4} \gamma \frac{1}{1-\gamma} + \frac{1}{4} \frac{1}{1-\gamma} = \frac{1}{4} \frac{1}{1-\gamma} (\gamma + 1) \tag{2.1}$$

$$= \frac{\gamma + 1}{4(1-\gamma)} \tag{2.2}$$

**Step 2 : Express  $V(g)$ ,  $V(0)$ , and  $V(1)$  in simplified forms**

$$V(g) = \zeta \sum_{i=0}^{\infty} \gamma^{2i} + \sum_{i=0}^{\infty} \gamma^{2i+1} \quad (2.3)$$

$$= \zeta \frac{1}{1-\gamma^2} + \gamma \frac{1}{1-\gamma^2} \quad (2.4)$$

$$= \frac{\zeta + \gamma}{1-\gamma^2} \quad (2.5)$$

$$V(0) = \zeta + \gamma^2 V(g) \quad (2.6)$$

$$= \zeta + \gamma^2 \frac{\zeta + \gamma}{1-\gamma^2} \quad (2.7)$$

$$= \frac{\zeta(1-\gamma^2) + \gamma^2(\zeta + \gamma)}{1-\gamma^2} \quad (2.8)$$

$$= \frac{\zeta + \gamma^3}{1-\gamma^2} \quad (2.9)$$

$$V(1) = \zeta + \gamma^2 V(0) \quad (2.10)$$

$$= \zeta + \gamma^2 \frac{\zeta + \gamma^3}{1-\gamma^2} \quad (2.11)$$

$$= \frac{\zeta(1-\gamma^2) + \gamma^2(\zeta + \gamma^3)}{1-\gamma^2} \quad (2.12)$$

$$= \frac{\zeta + \gamma^5}{1-\gamma^2} \quad (2.13)$$

**Step 3 : Substitute into the right side of the inequality**

$$\frac{1}{4}V(g) + \frac{2}{4}V(0) + \frac{1}{4}V(1) = \frac{1}{4} \frac{\zeta + \gamma}{1 - \gamma^2} + \frac{1}{2} \frac{\zeta + \gamma^3}{1 - \gamma^2} + \frac{1}{4} \frac{\zeta + \gamma^5}{1 - \gamma^2} \quad (2.14)$$

$$= \frac{1}{4(1 - \gamma^2)} [(\zeta + \gamma) + 2(\zeta + \gamma^3) + (\zeta + \gamma^5)] \quad (2.15)$$

$$= \frac{4\zeta + \gamma + 2\gamma^3 + \gamma^5}{4(1 - \gamma^2)} \quad (2.16)$$

**Step 4 : Set up the inequality**

$$\frac{\gamma + 1}{4(1 - \gamma)} \leq \frac{4\zeta + \gamma + 2\gamma^3 + \gamma^5}{4(1 - \gamma^2)} \quad (2.17)$$

**Step 5 : Use the identity  $1 - \gamma^2 = (1 - \gamma)(1 + \gamma)$** 

$$\frac{\gamma + 1}{4(1 - \gamma)} \leq \frac{4\zeta + \gamma + 2\gamma^3 + \gamma^5}{4(1 - \gamma)(1 + \gamma)} \quad (2.18)$$

**Step 6 : Multiply both sides by  $4(1 - \gamma)$** 

$$\gamma + 1 \leq \frac{4\zeta + \gamma + 2\gamma^3 + \gamma^5}{1 + \gamma} \quad (2.19)$$

**Step 7 : Multiply both sides by  $(1 + \gamma)$** 

$$(\gamma + 1)(1 + \gamma) \leq 4\zeta + \gamma + 2\gamma^3 + \gamma^5 \quad (2.20)$$

$$(\gamma + 1)^2 \leq 4\zeta + \gamma + 2\gamma^3 + \gamma^5 \quad (2.21)$$

**Step 8 : Expand and rearrange**

$$\gamma^2 + 2\gamma + 1 \leq 4\zeta + \gamma + 2\gamma^3 + \gamma^5 \quad (2.22)$$

$$4\zeta \geq \gamma^2 + 2\gamma + 1 - \gamma - 2\gamma^3 - \gamma^5 \quad (2.23)$$

$$4\zeta \geq \gamma^2 + \gamma + 1 - 2\gamma^3 - \gamma^5 \quad (2.24)$$

$$\zeta \geq \frac{\gamma^2 + \gamma + 1 - 2\gamma^3 - \gamma^5}{4} \quad (2.25)$$

Therefore, we obtain a **lower bound** on  $\zeta$  :

$$\zeta \geq \frac{\gamma^2 + \gamma + 1 - 2\gamma^3 - \gamma^5}{4} \quad (2.26)$$

where  $0 < \gamma < 1$ .

### 2.1.2 Step-by-step derivation of the upper bound on $\zeta$

Starting from the inequality :

$$\frac{1}{4}V(g) + \frac{1}{4}(\zeta + \gamma V(0)) + \frac{1}{4}(\zeta + \gamma V(1)) + \frac{1}{4}V(0) \leq \frac{1}{4}V(g) + \frac{2}{4}V(0) + \frac{1}{4}V(1)$$

**Step 1 : Cancel the  $\frac{1}{4}V(g)$  terms from both sides**

$$\frac{1}{4}(\zeta + \gamma V(0)) + \frac{1}{4}(\zeta + \gamma V(1)) + \frac{1}{4}V(0) \leq \frac{2}{4}V(0) + \frac{1}{4}V(1) \quad (2.27)$$

**Step 2 : Expand the left side**

$$\frac{1}{4}\zeta + \frac{1}{4}\gamma V(0) + \frac{1}{4}\zeta + \frac{1}{4}\gamma V(1) + \frac{1}{4}V(0) \leq \frac{2}{4}V(0) + \frac{1}{4}V(1) \quad (2.28)$$

**Step 3 : Combine like terms**

$$\frac{1}{2}\zeta + \frac{1}{4}\gamma V(0) + \frac{1}{4}\gamma V(1) \leq \frac{2}{4}V(0) - \frac{1}{4}V(0) + \frac{1}{4}V(1) \quad (2.29)$$

$$\frac{1}{2}\zeta + \frac{1}{4}\gamma V(0) + \frac{1}{4}\gamma V(1) \leq \frac{1}{4}V(0) + \frac{1}{4}V(1) \quad (2.30)$$

**Step 4 : Factor out common terms**

$$\frac{1}{2}\zeta \leq \frac{1}{4}V(0) + \frac{1}{4}V(1) - \frac{1}{4}\gamma V(0) - \frac{1}{4}\gamma V(1) \quad (2.31)$$

$$\frac{1}{2}\zeta \leq \frac{1}{4}V(0)(1 - \gamma) + \frac{1}{4}V(1)(1 - \gamma) \quad (2.32)$$

$$\frac{1}{2}\zeta \leq \frac{1 - \gamma}{4}(V(0) + V(1)) \quad (2.33)$$

$$\zeta \leq \frac{1 - \gamma}{2}(V(0) + V(1)) \quad (2.34)$$

**Step 5 : Substitute the expressions for  $V(0)$  and  $V(1)$**

$$V(0) + V(1) = \frac{\zeta + \gamma^3}{1 - \gamma^2} + \frac{\zeta + \gamma^5}{1 - \gamma^2} \quad (2.35)$$

$$= \frac{2\zeta + \gamma^3 + \gamma^5}{1 - \gamma^2} \quad (2.36)$$

**Step 6 : Substitute back into the inequality**

$$\zeta \leq \frac{1 - \gamma}{2} \cdot \frac{2\zeta + \gamma^3 + \gamma^5}{1 - \gamma^2} \quad (2.37)$$

$$= \frac{(1 - \gamma)(2\zeta + \gamma^3 + \gamma^5)}{2(1 - \gamma^2)} \quad (2.38)$$

**Step 7 : Use the identity  $1 - \gamma^2 = (1 - \gamma)(1 + \gamma)$**

$$\zeta \leq \frac{(1 - \gamma)(2\zeta + \gamma^3 + \gamma^5)}{2(1 - \gamma)(1 + \gamma)} \quad (2.39)$$

$$= \frac{2\zeta + \gamma^3 + \gamma^5}{2(1 + \gamma)} \quad (2.40)$$

**Step 8 : Multiply both sides by  $2(1 + \gamma)$**

$$2(1 + \gamma)\zeta \leq 2\zeta + \gamma^3 + \gamma^5 \quad (2.41)$$

$$2\zeta + 2\gamma\zeta \leq 2\zeta + \gamma^3 + \gamma^5 \quad (2.42)$$

$$2\gamma\zeta \leq \gamma^3 + \gamma^5 \quad (2.43)$$

$$\zeta \leq \frac{\gamma^3 + \gamma^5}{2\gamma} \quad (2.44)$$

$$\zeta \leq \frac{\gamma^2 + \gamma^4}{2} \quad (2.45)$$

Therefore, we obtain an **upper bound** on  $\zeta$  :

$$\zeta \leq \frac{\gamma^2 + \gamma^4}{2} \quad (2.46)$$

**Combined bounds :**

$$\frac{\gamma^2 + \gamma + 1 - 2\gamma^3 - \gamma^5}{4} \leq \zeta \leq \frac{\gamma^2 + \gamma^4}{2} \quad (2.47)$$

where  $0 < \gamma < 1$ .

### 2.1.3 Step-by-step derivation for the third inequality

Starting from the inequality :

$$\zeta \cdot \sum_{i=0}^{\infty} \gamma^i \leq \frac{1}{4}V(g) + \frac{2}{4}V(0) + \frac{1}{4}V(1)$$

**Step 1 : Simplify the left side using the geometric series**

$$\zeta \cdot \sum_{i=0}^{\infty} \gamma^i = \zeta \cdot \frac{1}{1-\gamma} \quad (2.48)$$

$$= \frac{\zeta}{1-\gamma} \quad (2.49)$$

**Step 2 : Use the previously derived expression for the right side** From our earlier calculation :

$$\frac{1}{4}V(g) + \frac{2}{4}V(0) + \frac{1}{4}V(1) = \frac{4\zeta + \gamma + 2\gamma^3 + \gamma^5}{4(1-\gamma^2)} \quad (2.50)$$

**Step 3 : Set up the inequality**

$$\frac{\zeta}{1-\gamma} \leq \frac{4\zeta + \gamma + 2\gamma^3 + \gamma^5}{4(1-\gamma^2)} \quad (2.51)$$

**Step 4 : Use the identity  $1 - \gamma^2 = (1 - \gamma)(1 + \gamma)$**

$$\frac{\zeta}{1-\gamma} \leq \frac{4\zeta + \gamma + 2\gamma^3 + \gamma^5}{4(1-\gamma)(1+\gamma)} \quad (2.52)$$

**Step 5 : Multiply both sides by  $(1 - \gamma)$**

$$\zeta \leq \frac{4\zeta + \gamma + 2\gamma^3 + \gamma^5}{4(1 + \gamma)} \quad (2.53)$$

**Step 6 : Multiply both sides by  $4(1 + \gamma)$**

$$4(1 + \gamma)\zeta \leq 4\zeta + \gamma + 2\gamma^3 + \gamma^5 \quad (2.54)$$

$$4\zeta + 4\gamma\zeta \leq 4\zeta + \gamma + 2\gamma^3 + \gamma^5 \quad (2.55)$$

**Step 7 : Subtract  $4\zeta$  from both sides**

$$4\gamma\zeta \leq \gamma + 2\gamma^3 + \gamma^5 \quad (2.56)$$

$$\zeta \leq \frac{\gamma + 2\gamma^3 + \gamma^5}{4\gamma} \quad (2.57)$$

$$\zeta \leq \frac{1 + 2\gamma^2 + \gamma^4}{4} \quad (2.58)$$

Therefore, we obtain another **upper bound** on  $\zeta$  :

$$\zeta \leq \frac{1 + 2\gamma^2 + \gamma^4}{4} \quad (2.59)$$

**Final combined bounds from all three inequalities :**

$$\frac{\gamma^2 + \gamma + 1 - 2\gamma^3 - \gamma^5}{4} \leq \zeta \leq \min \left\{ \frac{\gamma^2 + \gamma^4}{2}, \frac{1 + 2\gamma^2 + \gamma^4}{4} \right\} \quad (2.60)$$

where  $0 < \gamma < 1$ .

#### 2.1.4 Step-by-step derivation for the fourth inequality

Starting from the inequality :

$$\frac{1}{4}V(g) + \frac{1}{4}V(0) + \frac{1}{4}(\zeta + \gamma^2\zeta \sum_{i=0}^{\infty} \gamma^{2i}) + \frac{1}{4}(\zeta \sum_{i=0}^{\infty} \gamma^{2i}) \leq \frac{1}{4}V(g) + \frac{2}{4}V(0) + \frac{1}{4}V(1)$$

**Step 1 : Cancel the  $\frac{1}{4}V(g)$  terms from both sides**

$$\frac{1}{4}V(0) + \frac{1}{4}(\zeta + \gamma^2\zeta \sum_{i=0}^{\infty} \gamma^{2i}) + \frac{1}{4}(\zeta \sum_{i=0}^{\infty} \gamma^{2i}) \leq \frac{2}{4}V(0) + \frac{1}{4}V(1) \quad (2.61)$$

**Step 2 : Simplify using the geometric series  $\sum_{i=0}^{\infty} \gamma^{2i} = \frac{1}{1-\gamma^2}$**

$$\frac{1}{4}V(0) + \frac{1}{4}\left(\zeta + \gamma^2\zeta \frac{1}{1-\gamma^2}\right) + \frac{1}{4}\left(\zeta \frac{1}{1-\gamma^2}\right) \leq \frac{2}{4}V(0) + \frac{1}{4}V(1) \quad (2.62)$$

**Step 3 : Factor out common terms**

$$\frac{1}{4}V(0) + \frac{1}{4}\zeta\left(1 + \frac{\gamma^2}{1-\gamma^2}\right) + \frac{1}{4}\zeta \frac{1}{1-\gamma^2} \leq \frac{2}{4}V(0) + \frac{1}{4}V(1) \quad (2.63)$$

**Step 4 : Simplify the coefficient of  $\zeta$**

$$1 + \frac{\gamma^2}{1-\gamma^2} + \frac{1}{1-\gamma^2} = \frac{1-\gamma^2+\gamma^2+1}{1-\gamma^2} = \frac{2}{1-\gamma^2} \quad (2.64)$$

So the inequality becomes :

$$\frac{1}{4}V(0) + \frac{1}{4}\zeta \frac{2}{1-\gamma^2} \leq \frac{2}{4}V(0) + \frac{1}{4}V(1) \quad (2.65)$$

$$\frac{1}{4}V(0) + \frac{\zeta}{2(1-\gamma^2)} \leq \frac{1}{2}V(0) + \frac{1}{4}V(1) \quad (2.66)$$

**Step 5 : Rearrange to isolate the  $\zeta$  term**

$$\frac{\zeta}{2(1-\gamma^2)} \leq \frac{1}{2}V(0) - \frac{1}{4}V(0) + \frac{1}{4}V(1) \quad (2.67)$$

$$\frac{\zeta}{2(1-\gamma^2)} \leq \frac{1}{4}V(0) + \frac{1}{4}V(1) \quad (2.68)$$

$$\zeta \leq \frac{(1-\gamma^2)}{2}(V(0) + V(1)) \quad (2.69)$$



**Step 6 : Substitute the expressions for  $V(0)$  and  $V(1)$** 

$$V(0) + V(1) = \frac{\zeta + \gamma^3}{1 - \gamma^2} + \frac{\zeta + \gamma^5}{1 - \gamma^2} \quad (2.70)$$

$$= \frac{2\zeta + \gamma^3 + \gamma^5}{1 - \gamma^2} \quad (2.71)$$

**Step 7 : Substitute back into the inequality**

$$\zeta \leq \frac{(1 - \gamma^2)}{2} \cdot \frac{2\zeta + \gamma^3 + \gamma^5}{1 - \gamma^2} \quad (2.72)$$

$$= \frac{2\zeta + \gamma^3 + \gamma^5}{2} \quad (2.73)$$

**Step 8 : Multiply both sides by 2**

$$2\zeta \leq 2\zeta + \gamma^3 + \gamma^5 \quad (2.74)$$

$$0 \leq \gamma^3 + \gamma^5 \quad (2.75)$$

$$0 \leq \gamma^3(1 + \gamma^2) \quad (2.76)$$

Since  $0 < \gamma < 1$ , we have  $\gamma^3 > 0$  and  $(1 + \gamma^2) > 0$ , so this inequality is always satisfied. This means the fourth inequality does not provide an additional constraint on  $\zeta$ .

**Updated final bounds from all four inequalities :**

$$\frac{\gamma^2 + \gamma + 1 - 2\gamma^3 - \gamma^5}{4} \leq \zeta \leq \min \left\{ \frac{\gamma^2 + \gamma^4}{2}, \frac{1 + 2\gamma^2 + \gamma^4}{4} \right\} \quad (2.77)$$

where  $0 < \gamma < 1$ . The fourth inequality is automatically satisfied and does not further constrain the bounds.

## **2.2 Q-Learning**

## **2.3 Preferences over Decision Tree Policies**

## **2.4 Results**

## Conclusion

### **3.1 What happens when the MDP's transitions are independent of the current state?**



## **Deuxième partie**

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



# Chapitre 4

DPDT-intro





# Chapitre 5

DPDT-paper



# Chapitre 6

## Conclusion



## **Troisième partie**

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



## Conclusion générale





## Programmes informatiques

Les listings suivants sont au cœur de notre travail.

Listing A.1 – Il est l’heure

```

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

```

```
25 | printf("\n");  
26 | }
```

## Listing A.2 – Factorielle

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

# Table des matières

<b>Résumé</b>	<b>vii</b>
<b>Sommaire</b>	<b>ix</b>
<b>Preliminary Concepts</b>	<b>1</b>
Interpretable Sequential Decision Making . . . . .	1
What is Sequential Decision Making? . . . . .	1
What is Interpretability? . . . . .	2
What are existing approaches for learning interpretable programs?	4
Outline of the Thesis . . . . .	6
Technical Preliminaries . . . . .	6
What are decision trees? . . . . .	6
How to learn decision trees? . . . . .	7
Markov decision processes/problems . . . . .	7
Exact solutions for Markov decision problems . . . . .	9
Reinforcement learning of approximate solutions to MDPs . . . .	10
 <b>I A difficult problem : Learning Decision Trees for MDP</b>	 <b>11</b>
<b>1 A Decision Tree Policy for an MDP is a Policy for some Partially Observable MDP</b>	<b>13</b>
1.1 How to Learn a Decision Tree Policy for an MDP? . . . . .	14
1.1.1 Imitation . . . . .	14
1.1.2 Soft Trees . . . . .	14
1.1.3 Iterative Bounding MDPs . . . . .	14
1.2 How to solve Iterative Bounding MDPs? . . . . .	14
1.2.1 Asymmetric Reinforcement Learning . . . . .	14
1.2.2 Learning a decision tree policy is solving a POMDP . . . .	14
1.3 Is it hard to properly learn a Decision Tree Policy for an MDP? .	14
1.3.1 POMDPs are way harder to solve than MDPs . . . . .	14

1.3.2 Memoryless approaches to solve POMDPs seem ineffective	14
<b>2 An attempt at Learning Decision Tree Policies with Reinforcement Learning</b>	<b>15</b>
2.1 Grid Worlds . . . . .	15
2.1.1 Step-by-step derivation of the lower bound on $\zeta$ . . . . .	15
2.1.2 Step-by-step derivation of the upper bound on $\zeta$ . . . . .	18
2.1.3 Step-by-step derivation for the third inequality . . . . .	20
2.1.4 Step-by-step derivation for the fourth inequality . . . . .	21
2.2 Q-Learning . . . . .	24
2.3 Preferences over Decision Tree Policies . . . . .	24
2.4 Results . . . . .	24
<b>3 Conclusion</b>	<b>25</b>
3.1 What happens when the MDP's transitions are independent of the current state? . . . . .	25
 <b>II An easier problem : Learning Decision Trees for MDPs that are Classification tasks</b>	 <b>27</b>
4 DPDT-intro	29
5 DPDT-paper	31
6 Conclusion	33
 <b>III Beyond Decision Trees : what can be done with other Interpretable Policies?</b>	 <b>35</b>
Conclusion générale	37
A Programmes informatiques	39
Table des matières	41