



POLITECNICO DI MILANO  
DEPARTMENT OF INFORMATION, ELECTRONICS AND BIOENGINEERING  
DOCTORAL PROGRAMME IN INFORMATION TECHNOLOGY

---

ON THE EXPLOITATION OF UNCERTAINTY TO IMPROVE  
MAXIMUM EXPECTED VALUE ESTIMATE AND  
EXPLORATION IN REINFORCEMENT LEARNING

Doctoral Dissertation of:  
**Carlo D'Eramo**

Supervisor:  
**Prof. Marcello Restelli**

Tutor:  
**Prof. Andrea Bonarini**

The Chair of the Doctoral Program:  
**Prof. Marcello Restelli**

2019 – XXXI cycle



---

---

## Abstract

---

**A** BSTRACT goes here.



---

---

## Summary

---

SUMMARY goes here.



---

---

## Contents

---

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Perception and interaction . . . . .	1
1.2	Learn how to act with Reinforcement Learning . . . . .	2
1.2.1	Uncertainty in Reinforcement Learning . . . . .	2
1.2.2	Balancing exploration and exploitation . . . . .	3
1.3	My research . . . . .	3
<b>2</b>	<b>Preliminaries</b>	<b>5</b>
<b>3</b>	<b>Maximum Expected Value estimate</b>	<b>7</b>
<b>4</b>	<b>Exploration</b>	<b>9</b>
<b>5</b>	<b>Deep</b>	<b>11</b>
<b>6</b>	<b>Mushroom</b>	<b>13</b>
<b>7</b>	<b>Conclusion</b>	<b>15</b>
	<b>Bibliography</b>	<b>17</b>





---

---

## List of Figures

---

1.1 Reinforcement Learning problem scheme . . . . .	2
---	---



---

---

## List of Tables

---



---

# CHAPTER 1

---

## Introduction

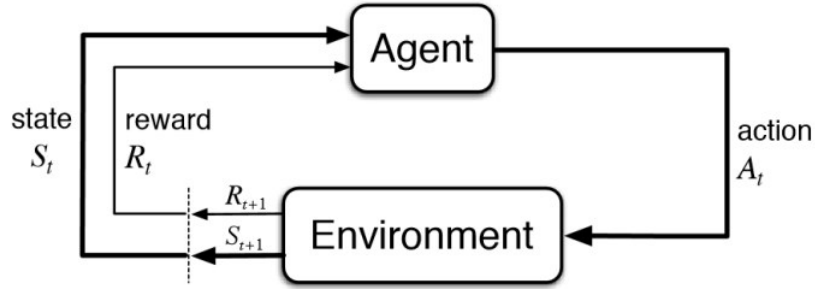
---

EVERYONE experiences the process of taking decisions during his life. As a matter of fact, drastically the life of an individual can be synthesized in its *perception* of the world and its *interaction* with it. The concepts of perception and interaction might seem quite straightforward to understand: for a human the perception of the world comes from its senses and the interaction comes from its possibility to change its surroundings. On the contrary, these concepts are actually absolutely hard to define and aroused, during the centuries, a strong debate between scientists, biologists, and even philosophers.

### 1.1 Perception and interaction

---

We start from the assumption that, by definition, an individual perceive the environment around it and acts on it in order to achieve *goals* expressed by its will. In other words, all the actions made by an individual are done to satisfy its will to obtain something from the world it lives in. This task is naturally performed by humans, but it implies some challenging problems that are hard, or unfeasible, to solve. One of them comes from the intrinsic *uncertainty* of the perception we have of the world around us. Indeed, the perception of the world consists in the interpretation of the information provided by senses, but the process of information retrieval by senses and the mental processes to understand them, inevitably introduce a certain level of noise that distorts the original true information. On the other hand, the interaction with the world deals with the will of the individual to perform actions to change the environment around it, but this apparently simple operation involves complex biologic mechanisms to coordinate the body according to the will and difficulties in the perception of the consequences of the interaction. Moreover, the concept of goal can be unclear and the individual may



**Figure 1.1:** *The scheme of a RL model.*

result in performing actions without being sure of what it want. It is arguable that discussing about the concept of true information and the concept of will requires strong theoretical considerations since they are both hardly definable concepts. For many centuries scientists and philosophers debated about these topics, in particular trying to solve complex problems like the real nature of perceivable things and the concept of free will. However, to make the discussion of these concepts suitable for our purposes throughout all this thesis, we lighten the definition of them to the one provided by common sense.

## 1.2 Learn how to act with Reinforcement Learning

---

Reinforcement Learning (RL) [5] is a subfield of Machine Learning (ML) which aims to realize autonomous *agents* able to learn how to act in a certain *environment* in order to maximize an objective function; this is achieved providing the agent with the perception of its *state* in the environment and making it learn the appropriate *action* to perform, where an action can be seen as an atomic operation that brings the agent from a state to another one. The objective function represents a measure of how well the agent is accomplishing its task in the environment and it is usually formalized by a discounted sum of *rewards* obtained after each action (Figure 1.1). The sum is discounted to give more importance to the most recent rewards w.r.t. the ones further in the future. The reward function, i.e. the function returning the reward after each action, is not a concrete part of the environment, but it is often designed by a human which decides whether a certain behavior has to be reinforced (returning a positive reward) or inhibited (returning a negative reward).

### 1.2.1 Uncertainty in Reinforcement Learning

The major challenge of RL is represented by the uncertainty. In fact, initially, the agent is not provided with the knowledge of how the environment will react to its action, thus it does not know whether an action would be good or not to maximize its objective function. In other words, before trying an action, it does not know if that action will get a positive or a negative reward, and it does not know if that action will let it go to the desired state or not. Thus, the former problem can be seen as uncertainty in the reward function and the latter as uncertainty in the transition (i.e. model) function. In some cases, also the uncertainty in the perception of the current state of the agent is

considered, making the problem more complex.

The uncertainty issue results in the need of the agent to try actions in order to improve its knowledge of the environment. This process delays the collection of high rewards, but helps the agent to reduce its uncertainty. However, since the objective function is a sum of discounted rewards where later rewards worth less than recent ones, the agent also needs to learn fast in order to learn to perform the most rewarding actions as soon as possible. The need to explore and the need to *exploit* the actions believed to be good introduces an important problem known as *exploration-exploitation dilemma*.

### **1.2.2 Balancing exploration and exploitation**

The exploration-exploitation dilemma has been broadly studied in the field of Multi-Armed Bandit (MAB), a particular case of the RL problem with a single state [4]. In this problem the goal is to find the sequence of optimal actions, i.e. the sequence of actions that allows to maximize the return. The simplistic setting of the MAB problem allows to theoretically study the balancing of exploratory and exploitative actions, for instance to derive upper confidence bounds on the *regret*, i.e. a measure of the return lost in performing non-optimal actions [1, 3, 7], and several algorithms to address this problem have been proposed such as UCB1 [2] and Thompson sampling [6].

The RL setting complicates the MAB problem because of the presence of multiple states. This makes the exploration-exploitation dilemma less tractable in terms of complexity and computational feasibility. Indeed, the quality of the actions must now be evaluated for each state, contrarily to the MAB case where the presence of a single state simplifies the problem. This issue is what makes RL so challenging and has been addressed for decades in the literature by lots of methodologies and algorithms.

## **1.3 My research**

---





---

## CHAPTER 2

---

### Preliminaries

---



---

## CHAPTER 3

---

### **Maximum Expected Value estimate**

---



---

# CHAPTER 4

---

## Exploration

---



---

# CHAPTER 5

---

**Deep**

---





---

## CHAPTER 6

---

### Mushroom

---



---

# CHAPTER 7

---

## Conclusion

---



---

---

## Bibliography

---

- [1] Shipra Agrawal and Navin Goyal. Analysis of thompson sampling for the multi-armed bandit problem. In *Conference on Learning Theory*, pages 39–1, 2012.
- [2] Peter Auer, Nicolo Cesa-Bianchi, and Paul Fischer. Finite-time analysis of the multiarmed bandit problem. *Machine learning*, 47(2-3):235–256, 2002.
- [3] Sébastien Bubeck, Nicolo Cesa-Bianchi, et al. Regret analysis of stochastic and nonstochastic multi-armed bandit problems. *Foundations and Trends® in Machine Learning*, 5(1):1–122, 2012.
- [4] Tze Leung Lai and Herbert Robbins. Asymptotically efficient adaptive allocation rules. *Advances in applied mathematics*, 6(1):4–22, 1985.
- [5] Richard S Sutton, Andrew G Barto, et al. *Reinforcement learning: An introduction*. MIT press, 1998.
- [6] William R Thompson. On the likelihood that one unknown probability exceeds another in view of the evidence of two samples. *Biometrika*, 25(3/4):285–294, 1933.
- [7] Joannes Vermorel and Mehryar Mohri. Multi-armed bandit algorithms and empirical evaluation. In *European conference on machine learning*, pages 437–448. Springer, 2005.