

# Optimal sequential decision making for complex problems

Damien Ernst – University of Liège  
Email: [dernst@uliege.be](mailto:dernst@uliege.be)

## About the class

Regular lectures notes about various topics on the subject with a great emphasis on artificial intelligence and the design of **intelligent agents**.

Exercises (bring your computer) and projects to illustrate the concepts seen during the class.

Discussing together of research papers and going through presentations of other authors

## Artificial autonomous intelligent agent: formal definition

An **agent** is anything that is capable of acting upon information it perceives.

An **intelligent agent** is an agent capable of making decisions about how it acts based on experience, that is of learning decision from experience.

An **autonomous intelligent agent** is an intelligent agent that is free to choose between different actions.

An **artificial autonomous intelligent agent** is anything we create that is capable of actions based on information it perceives, its own experience, and its own decisions about which actions to perform.

Since “artificial autonomous intelligent agent” is quite mouthful, we follow the convention of using “intelligent agent” or “autonomous agent” for short.

## Application of intelligent agents

**Intelligent agents** are applied in variety of areas: project management, electronic commerce, robotics, information retrieval, military, networking, planning and scheduling, etc.

Examples:

- A **web search agent** that may have the goal of obtaining web site addresses that would match the query of history made by customer. It could operate in the background and deliver recommendations to the customer on a weekly basis.
- A **robot agent** that would learn to fulfill some specific tasks through experience such as playing soccer, cleaning, etc.
- An intelligent **flight control** system
- An agent **for placing advertisements on the web**.
- A **intelligent agent for playing computer games** that does rely on some predefined strategies for playing but that would learn them by interacting with some opponents.

## Machine learning and reinforcement learning: definitions

**Machine learning** is a broad subfield of artificial intelligence is concerned with the development of algorithms and techniques that allow computers to "learn".

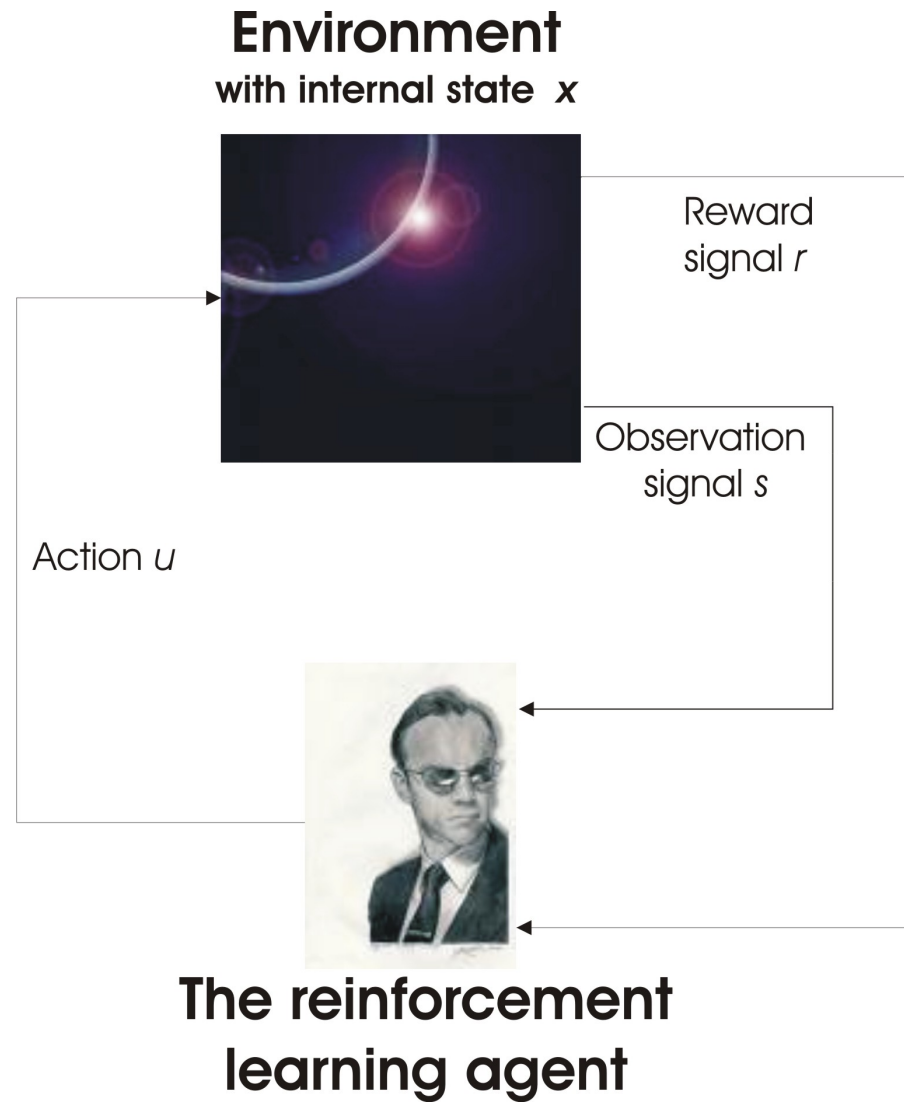
**Reinforcement Learning** (RL in short) refers to a class of problems in machine learning which postulate an autonomous agent exploring an environment in which the agent perceives information about its current state and takes actions. The environment, in return, provides a reward signal (which can be positive or negative). The agent has as objective to maximize the (expected) cumulative reward signal over the course of the interaction.

The **policy** of an agent determines the way the agent selects its action based on the information it has. A policy can be either deterministic or stochastic.

**Research in reinforcement learning** aims at designing policies which lead to large (expected) cumulative reward.

**Where does the intelligence come from ?** The policies process in an “**intelligent way**” the information to select “**good actions**”.

# An RL agent interacting with its environment



## Some generic difficulties with designing intelligent agents

- **Inference problem.** The environment **dynamics** and the mechanism behind the **reward** signal are (partially) **unknown**. The policies need to be able to infer from the information the agent has gathered from interaction with the system, “good control actions”.
- **Computational complexity.** The policy must be able to process the history of the observation within limited amount of computing times and memory.
- **Tradeoff between exploration and exploitation.\*** To obtain a lot of reward, a reinforcement learning agent must prefer actions that it has tried in the past and found to be effective in producing reward. But to discover such actions, it has to try actions that it has not selected before.

\*May be seen as a subproblem of the general inference problem. This problem is often referred to in the “classical control theory” as the dual control problem.



The agent has to exploit what it already knows in order to obtain reward, but it also has to explore in order to make better action selections in the future. The dilemma is that neither exploration nor exploitation can be pursued exclusively without failing at the task. The agent must try a variety of actions and progressively favor those that appear to be best. On a stochastic task, each action must be tried many times to gain a reliable estimate its expected reward.

- **Exploring safely the environment.** During an exploration phase (more generally, any phase of the agent's interaction with its environment), the agent must avoid reaching unacceptable states (e.g., states that may for example endanger its own integrity). By associating rewards of  $-\infty$  to those states, exploring safely can be assimilated to a problem of exploration-exploitation.
- **Adversarial environment.** The environment may be adversarial. In such a context, one or several other players seek to adopt strategies that oppose the interests of the RL agent.

## Different characterizations of RL problems

- **Stochastic** (e.g.,  $x_{t+1} = f(x_t, u_t, w_t)$  where the random disturbance  $w_t$  is drawn according to the conditional probability distribution  $P_w(\cdot|x_t, u_t)$ ) versus **deterministic** (e.g.,  $x_{t+1} = f(x_t, u_t)$ )
- **Partial observability** versus **full observability**. The environment is said to be partially (fully) observable if the signal  $s_t$  describes partially (fully) the environment's state  $x_t$  at time  $t$ .
- **Time-invariant** (e.g.,  $x_{t+1} = f(x_t, u_t, w_t)$  with  $w_t = P_w(\cdot|x_t, u_t)$ ) versus **time-variant** (e.g.,  $x_{t+1} = f(x_t, u_t, w_t, t)$ ) dynamics.
- **Continuous** (e.g.,  $\dot{x} = f(x, u, w)$ ) versus **discrete** dynamics (e.g.,  $x_{t+1} = f(x_t, u_t, w_t)$ ).

- **Multi-agent** framework versus **single-agent** framework. In a multi-agent framework the environment may be itself composed of (intelligent) agents. A multi-agent framework can often be assimilated to a single-agent framework by considering that the internal states of the other agents are unobservable variables. **Game theory** and, more particularly, the theory of **learning in games** study situations where various intelligent agents interact with each other.
- **Finite time** versus **infinite time** of interaction.
- **Single state** versus **multi-state** environment. In single state environment, computation of an optimal policy for the agent is often reduced to the computation of the maximum of a stochastic function (e.g., find  $u^* \in \arg \max_{u \in U} E_{w \sim P_w(\cdot|u)} [r(u, w)]$ ).

- **Multi-objective** reinforcement learning agent (reinforcement learning signal can be multi-dimensional) versus **single-objective** RL agent.
- **Risk-adverse** reinforcement learning agent. The goal of the agent is not anymore to maximize the expected cumulative reward but maximize the lowest cumulative reward it could possibly obtain.

## Characterization of the RL problem adopted in this class

- **Dynamics of the environment:**

$$x_{t+1} = f(x_t, u_t, w_t) \quad t = 0, 1, 2 \dots$$

where for all  $t$ , the state  $x_t$  is an element of the state space  $X$ , the action  $u_t$  is an element of the action space  $U$  and the random disturbance  $w_t$  is an element of the disturbance space  $W$ .

Disturbance  $w_t$  generated by the time-invariant conditional probability distribution  $P_w(\cdot|x, u)$ .

- **Reward signal:**

The function  $r(x, u, w)$  is the so-called **reward function** supposed to be bounded by below by 0 and by above by a constant  $B_r \geq 0$ .

To the transition from  $t$  to  $t + 1$  is associated a reward signal  $\gamma^t r_t = \gamma^t r(x_t, u_t, w_t)$  where  $r(x, u, w)$  is a reward function supposed to be bounded by a constant  $B_r$  and  $\gamma \in [0, 1[$  a decay factor.

- **Cumulative reward signal:**

Let  $h_t \in \mathcal{H}$  be the trajectory from instant time 0 to  $t$  in the combined state, action, reward spaces:

$h_t = (x_0, u_0, r_0, x_1, u_1, r_1, \dots, u_{t-1}, r_{t-1}, x_t)$ . Let  $\pi \in \Pi$  be a stochastic policy such that  $\pi : \mathcal{H} \times U \rightarrow [0, 1]$  and let us denote by  $J^\pi(x)$  the expected return of a policy  $\pi$  (or **expected cumulative reward signal**) when the system starts from  $x_0 = x$

$$J^\pi(x) = \lim_{T \rightarrow \infty} E\left[\sum_{t=0}^T \gamma^t r(x_t, u_t \sim \pi(h_t, \cdot), w_t) \mid x_0 = x\right]$$

- **Information available:**

The agent does not know  $f$ ,  $r$  and  $P_w$ . The only information it has on these three elements is the information contained in  $h_t$ .

## Goal of reinforcement learning

- Let  $\pi^* \in \Pi$  a policy such that  $\forall x \in X$ ,

$$J^{\pi^*}(x) = \max_{\pi \in \Pi} J^\pi(x) \quad (I)$$

Under some mild assumptions\* on  $f$ ,  $r$  and  $P_w$ , such a policy  $\pi^*$  indeed exists.

- In reinforcement learning, we want to build policies  $\hat{\pi}^*$  such that  $J^{\hat{\pi}^*}$  is as close as possible (according to specific metrics) to  $J^{\pi^*}$ .
- If  $f$ ,  $r$  and  $P_w$  were known, we could, by putting aside the difficulty of finding in  $\Pi$  the policy  $\pi^*$ , design the optimal agent by solving the optimal control problem (I). However,  $J^\pi$  depends on  $f$ ,  $r$  and  $P_w$  which are supposed to be unknown  $\Rightarrow$  How can we solve this combined inference - optimization problem ?

\*We will suppose that these mild assumptions are always satisfied afterwards.

## Dynamic Programming (DP) theory reminder: optimality of stationary policies

- A stationary control policy  $\mu : X \rightarrow U$  selects at time  $t$  the action  $u_t = \mu(x_t)$ . Let  $\Pi_\mu$  denote the set of stationary policies.
- The *expected return* of a stationary policy when the system starts from  $x_0 = x$  is

$$J^\mu(x) = \lim_{T \rightarrow \infty} E_{w_0, w_1, \dots, w_T} \left[ \sum_{t=0}^T \gamma^t r(x_t, \mu(x_t), w_t) \mid x_0 = x \right]$$

- Let  $\mu^*$  be a policy such that  $J^{\mu^*}(x) = \max_{\mu \in \Pi_\mu} J^\mu(x)$  everywhere on  $X$ .

It can be shown that such a policy indeed exists. We name such a policy an *optimal stationary policy*.

- From classical dynamic programming theory, we know  $J^{\mu^*}(x) = J^{\pi^*}(x)$  everywhere  $\Rightarrow$  considering only stationary policies is not suboptimal!



## From Dynamic Programming (DP): How to compute the return of a stationary policy

- We define the **functions**  $J_N^\mu : X \rightarrow \mathbb{R}$  by the recurrence equation

$$J_N^\mu(x) = \mathop{E}_{w \sim P_w(\cdot|x,u)} [r(x, \mu(x), w) + \gamma J_{N-1}^\mu(f(x, \mu(x), w))], \quad \forall N \geq 1 \quad (1)$$

with  $J_0^\mu(x) \equiv 0$ .

- We have as result of the DP theory

$$\lim_{N \rightarrow \infty} \|J_N^\mu - J^\mu\|_\infty \rightarrow 0 \quad (2)$$

and

$$\|J_N^\mu - J^\mu\|_\infty \leq \frac{\gamma^N}{1 - \gamma} B_r \quad (3)$$

## DP theory reminder: $Q_N$ -functions and $\mu^*$

- We define the **functions**  $Q_N : X \times U \rightarrow \mathbb{R}$  by the recurrence equation

$$Q_N(x, u) = \mathop{E}_{w \sim P_w(\cdot|x, u)} [r(x, u, w) + \gamma \max_{u' \in U} Q_{N-1}(f(x, u, w), u')], \quad \forall N \geq 1 \quad (4)$$

with  $Q_0(x, u) \equiv 0$ . These  $Q_N$ -functions are also known as **state-action value** functions.

- We denote by  $\mu_N^* : X \rightarrow U$  the stationary policy:

$$\mu_N^*(x) \in \arg \max_{u \in U} Q_N(x, u). \quad (5)$$

- We define the  $Q$ -function as being the unique solution of the **Bellman equation**:

$$Q(x, u) = \mathop{E}_{w \sim P_w(\cdot|x, u)} [r(x, u, w) + \gamma \max_{u' \in U} Q(f(x, u, w), u')]. \quad (6)$$

We have the following results:

- **Convergence** in infinite norm of the sequence of functions  $Q_N$  to the  $Q$ -function, i.e.  $\lim_{N \rightarrow \infty} \|Q_N - Q\|_\infty \rightarrow 0$  (see Appendix **I** for the proof)

- A control policy  $\mu^*$  is optimal **if and only if**

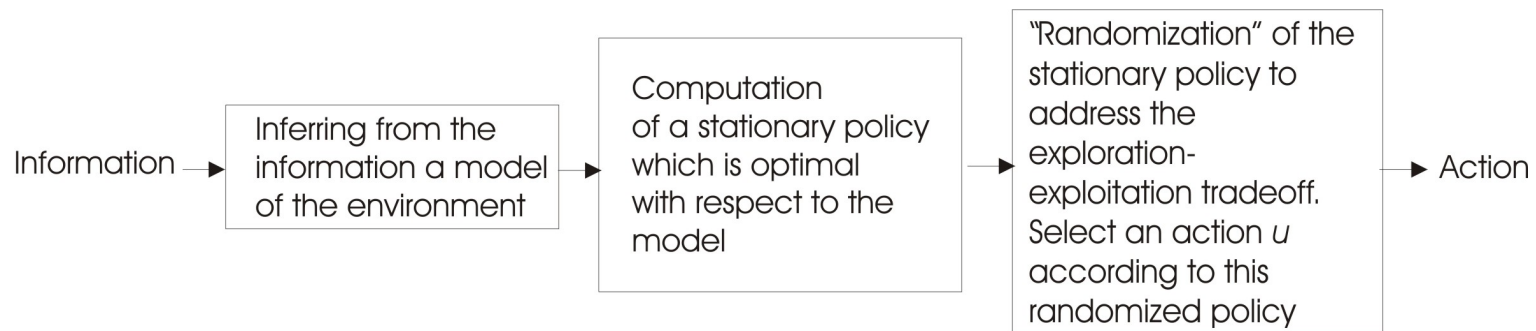
$$\mu^*(x) \in \arg \max_{u \in U} Q(x, u) \quad (7)$$

We also have  $J^{\mu^*}(x) = \max_{u \in U} Q(x, u)$ , which is also called the **value function**.

- The following **bound on the suboptimality of  $\mu_N^*$**  with respect to  $\mu^*$  holds:

$$\|J^{\mu_N^*} - J^{\mu^*}\|_\infty \leq \frac{2\gamma^N B_r}{(1 - \gamma)^2} \quad (8)$$

## A pragmatic approach for designing (hopefully) good policies $\hat{\pi}^*$



We focus first on to the design of functions  $\hat{\pi}^*$  which realize sequentially the following three tasks:

1. "System identification" phase. Estimation from  $h_t$  of an approximate system dynamics  $\hat{f}$ , an approximate probability distribution  $\hat{P}_w$  and an approximate reward function  $\hat{r}$ .

## 2. Resolution of the optimization problem:

Find in  $\Pi_\mu$  the policy  $\hat{\mu}^*$  such that  $\forall x \in X, J^{\hat{\mu}^*}(x) = \max_{\mu \in \Pi_\mu} \hat{J}^\mu(x)$

where  $\hat{J}^\mu$  is defined similarly as function  $J^\mu$  but with  $\hat{f}$ ,  $\hat{P}_w$  and  $\hat{r}$  replacing  $f$ ,  $P_w$  and  $r$ , respectively.

3. Afterwards, the policy  $\hat{\pi}$  selects with a probability  $1 - \epsilon(h_t)$  actions according to the policy  $\hat{\mu}^*$  and with a probability  $\epsilon(h_t)$  at random. Step 3 has been introduced to address the **dilemma between exploration and exploitation**.\*

\*We will not address further the design of the 'right function'  $\epsilon : \mathcal{H} \rightarrow [0, 1]$ . In many applications, it is chosen equal to a small constant (say, 0.05) everywhere.

## Some constructive algorithms for designing $\hat{\pi}^*$ when dealing with finite state-action spaces

- Until say otherwise, we consider the particular case of finite state and action spaces (i.e.,  $X \times U$  finite).
- When  $X$  and  $U$  are finite, there exists a vast panel of 'well-working' implementable RL algorithms.
- We focus first on approaches which solve separately Step **1.** and Step **2.** and then on approaches which solve both steps together.
- The proposed algorithms infer  $\hat{\mu}^*$  from  $h_t$ . They can be adapted in a straightforward way to episode-based reinforcement learning where a model of  $\mu^*$  must be inferred from several trajectories  $h_{t_1}, h_{t_2}, \dots, h_{t_m}$  with  $t_i \in \mathbb{N}_0$ .

## Reminder on Markov Decision Processes

- A Markov Decision Process (**MDP**) is defined through the following objects: a state space  $X$ , an action space  $U$ , transition probabilities  $p(x'|x, u) \forall x, x' \in X, u \in U$  and a reward function  $r(x, u)$ .
- $p(x'|x, u)$  gives the probability of reaching state  $x'$  after taking action  $u$  while being in state  $x$ .
- We consider MDPs for which we want to find decision policies that maximize the reward signal  $\gamma^t r(x_t, u_t)$  over an infinite time horizon.
- MDPs can be seen as a **particular type** of the discrete-time optimal control problem introduced earlier where the system dynamics is expressed under the form of transition probabilities and where the reward function does not depend on the disturbance  $w$  anymore.

## MDP structure definition from the system dynamics and reward function

- We define\*

$$r(x, u) = E_{w \sim P_w(\cdot|x, u)} [r(x, u, w)] \quad \forall x \in X, u \in U \quad (9)$$

$$p(x'|x, u) = E_{w \sim P_w(\cdot|x, u)} [I_{\{x'=f(x, u, w)\}}] \quad \forall x, x' \in X, u \in U \quad (10)$$

- Equations (9) and (10) define the **structure of an equivalent MDP** in the sense that the expected return of any policy applied to the original optimal control problem is equal to its expected return for the MDP.

- The **recurrence equation** defining the functions  $Q_N$  can be rewritten:

$$Q_N(x, u) = r(x, u) + \gamma \sum_{x' \in X} p(x'|x, u) \max_{u' \in U} Q_{N-1}(x', u'), \quad \forall N \geq 1 \text{ with}$$

$$Q_0(x, u) \equiv 0.$$

\* $I_{\{\text{logical\_expression}\}} = 1$  if *logical\_expression* is true and 0 if *logical\_expression* is false.



## Reminder: Random variable and strong law of large numbers

- A **random variable** is not a variable but rather a function that maps outcomes (of an experiment) to numbers. Mathematically, a random variable is defined as a measurable function from a **probability space** to some **measurable space**. We consider here random variables  $\theta$  defined on the probability space  $(\Omega, P)$ .\*

- $\frac{E}{P}[\theta]$  is the mean value of the random variable  $\theta$ .

- Let  $\theta_1, \theta_2, \dots, \theta_n$  be  $n$  values of the random variable  $\theta$  which are drawn independently. Suppose also that  $\frac{E}{P}[|\theta|] = \int_{\Omega} |\theta| dP$  is smaller than  $\infty$ . In such a case, the **strong law of large number** states that:

$$\lim_{n \rightarrow \infty} \frac{\theta_1 + \theta_2 + \dots + \theta_n}{n} \xrightarrow{P} \frac{E}{P}[\theta] \quad (11)$$

\*For the sake of simplicity, we have considered here that  $(\Omega, P)$  indeed defines a probability space which is not rigorous.

## Step 1. Identification by learning the structure of the equivalent MPD

- The objective is to infer some 'good approximations' of  $p(x'|x, u)$  and  $r(x, u)$  from:

$$h_t = (x_0, u_0, r_0, x_1, u_1, r_1, \dots, u_{t-1}, r_{t-1}, x_t)$$

### Estimation of $r(x, u)$ :

Let  $A(x, u) = \{k \in \{0, 1, \dots, t-1\} | (x_k, u_k) = (x, u)\}$ . Let  $k_1, k_2, \dots, k_{\#A(x, u)}$  denote the elements of the set.\* The values  $r_{k_1}, r_{k_2}, \dots, r_{k_{\#A(x, u)}}$  are  $\#A(x, u)$  values of the random variable  $r(x, u, w)$  which are drawn independently. It follows therefore naturally that to estimate its mean value  $r(x, u)$ , we can use the following unbiased estimator:

$$\hat{r}(x, u) = \frac{\sum_{k \in A(x, u)} r_k}{\#A(x, u)} \quad (12)$$

\*If  $S$  is a set of elements,  $\#S$  denote the cardinality of  $S$ .

Estimation of  $p(x'|x, u)$ :

The values  $I_{\{x'=x_{k_1+1}\}}$ ,  $I_{\{x'=x_{k_2+1}\}}$ ,  $\dots$ ,  $I_{\{x'=x_{k_{\#A(x,u)+1}\}}$  are  $\#A(x, u)$  values of the random variable  $I_{\{x'=f(x,u,w)\}}$  which are drawn independently. To estimate its mean value  $p(x'|x, u)$ , we can use the unbiased estimator:

$$\hat{p}(x'|x, u) = \frac{\sum_{k \in A(x, u)} I_{\{x_{k+1}=x'\}}}{\#A(x, u)} \quad (13)$$

## Step 2. Computation of $\hat{\mu}^*$ identification by learning the structure of the equivalent MPD

- We compute the  $\hat{Q}_N$ -functions from the knowledge of  $\hat{r}$  and  $\hat{p}$  by exploiting the **recurrence equation**:

$$\hat{Q}_N(x, u) = \hat{r}(x, u) + \gamma \sum_{x' \in X} \hat{p}(x'|x, u) \max_{u' \in U} Q_{N-1}(x', u'), \quad \forall N \geq 1 \text{ with}$$

$$\hat{Q}_0(x, u) \equiv 0 \text{ and then take}$$

$$\hat{\mu}_N^* = \arg \max_{u \in U} \hat{Q}_N(x, u) \quad \forall x \in X \quad (14)$$

as approximation of the optimal policy, with  $N$  'large enough' (e.g., right hand side of inequality (8) drops below  $\epsilon$ ).

- One can show that if the estimated MDP structure lies in an ' **$\epsilon$ -neighborhood**' of the true structure, then,  $J^{\hat{\mu}^*}$  is in a ' **$O(\epsilon)$ -neighborhood**' of  $J^{\mu^*}$  where  $\hat{\mu}^*(x) = \lim_{N \rightarrow \infty} \arg \max_{u \in U} \hat{Q}_N(x, u)$ .

## The case of limited computational resources

- **Number of operations** to estimate the MDP structure grows linearly with  $t$ . **Memory requirements** needed to store  $h_t$  also **grow linearly with  $t$**   $\Rightarrow$  an agent having limited computational resources will face problems after certain time of interaction.

- We describe an algorithm which requires at time  $t$  a number of operations that does not depend on  $t$  to update the MDP structure and for which the memory requirements do not grow with  $t$ :

**At time 0**, set  $N(x, u) = 0$ ,  $N(x, u, x') = 0$ ,  $R(x, u) = 0$ ,  $p(x'|x, u) = 0$ ,  $\forall x, x' \in X$  and  $u \in U$ .

**At time  $t \neq 0$** , do

1.  $N(x_{t-1}, u_{t-1}) \leftarrow N(x_{t-1}, u_{t-1}) + 1$
2.  $N(x_{t-1}, u_{t-1}, x_t) \leftarrow N(x_{t-1}, u_{t-1}, x_t) + 1$
3.  $R(x_{t-1}, u_{t-1}) \leftarrow R(x_{t-1}, u_{t-1}) + r_t$
4.  $r(x_{t-1}, u_{t-1}) \leftarrow \frac{R(x_{t-1}, u_{t-1})}{N(x_{t-1}, u_{t-1})}$
5.  $p(x|x_{t-1}, u_{t-1}) \leftarrow \frac{N(x_{t-1}, u_{t-1}, x)}{N(x_t, u_t)} \quad \forall x \in X$

## Merging Step 1. and 2. to learn directly the $Q$ -function: the $Q$ -learning algorithm

The  $Q$ -learning algorithm is an algorithm that infers directly from  $h_t = (x_0, u_0, r_0, x_1, u_1, r_1, \dots, u_{t-1}, r_{t-1}, x_t)$  an approximate value of the  $Q$ -function, without identifying the structure of a Markov Decision Process.

The algorithm can be described by the following steps:

1. Initialisation of  $\hat{Q}(x, u)$  to 0 everywhere. Set  $k = 0$ .
2.  $\hat{Q}(x_k, u_k) \leftarrow (1 - \alpha_k)\hat{Q}(x_k, u_k) + \alpha_k(r_k + \gamma \max_{u \in U} \hat{Q}_k(x_{k+1}, u))$
3.  $k \leftarrow k + 1$ . If  $k = t$ , return  $\hat{Q}$  and stop. Otherwise, go back to 2.

## Q-learning: some remarks

- Iteration 2. can be rewritten as  $\hat{Q}(x_k, u_k) \leftarrow \hat{Q}(x_k, u_k) + \alpha \delta(x_k, u_k)$  where the term:

$$\delta(x_k, u_k) = r_k + \gamma \max_{u \in U} \hat{Q}(x_{k+1}, u) - \hat{Q}(x_k, u_k), \quad (15)$$

called the *temporal difference*.

- **Learning ratio  $\alpha_k$** : The learning ratio  $\alpha_k$  is often chosen constant with  $k$  and equal to a small value (e.g.,  $\alpha_k = 0.05, \forall k$ ).
- **Consistency of the Q-learning algorithm**: Under some particular conditions on the way  $\alpha_k$  decreases to zero ( $\lim_{t \rightarrow \infty} \sum_{k=0}^{t-1} \alpha_k \rightarrow \infty$  and  $\lim_{t \rightarrow \infty} \sum_{k=0}^{t-1} \alpha_k^2 < \infty$ ) and the history  $h_t$  (when  $t \rightarrow \infty$ , every state-action pair needs to be visited an infinite number of times),  $\hat{Q} \rightarrow Q$  when  $t \rightarrow \infty$ .

- **Experience replay:** At each iteration, the  $Q$ -learning algorithm uses a sample  $l_k = (x_k, u_k, r_k, x_{k+1})$  to update the function  $\hat{Q}$ . If rather than to use the finite sequence of sample  $l_0, l_2, \dots, l_{t-1}$ , we use the infinite size sequence  $l_{i_1}, l_{i_2}, \dots$  to update in a similar way  $\hat{Q}$ , where the  $i_j$  are i.i.d. with uniform distribution on  $\{0, 2, \dots, t-1\}$ , then  $\hat{Q}$  converges to the approximate  $Q$ -function computed from the estimated equivalent MDP structure.



## Inferring $\hat{\mu}^*$ from $h_t$ when dealing with very large or infinite state-action spaces

- Up to now, we have considered problems having discrete (and not too large) state and action spaces  $\Rightarrow \hat{\mu}^*$  and the  $\hat{Q}_N$ -functions could be represented in a **tabular form**.
- We consider now the case of **very large or infinite state-action spaces**: functions approximators need to be used to represent  $\hat{\mu}^*$  and the  $\hat{Q}_N$ -functions.
- These function approximators need to be used in a way that there are able to '**well generalize**' over the whole state-action space the information contained in  $h_t$ .
- There is a vast literature on function approximators in reinforcement learning. We focus first on one algorithm named '**fitted Q iteration**' which computes the functions  $\hat{Q}_N$  from  $h_t$  by solving a sequence of batch mode supervised learning problems.

## Reminder: Batch mode supervised learning

- A batch mode Supervised Learning (SL) algorithm infers from a set of **input-output** (**input** = information state); (**output** = class label, real number, graph, etc) a model which explains “at best” these input-output pairs.
- **A loose formalisation of the SL problem:** Let  $I$  be the input space,  $O$  the output space,  $\Xi$  the disturbance space. Let  $g : I \times \Xi \rightarrow O$ . Let  $P_\xi(\cdot|i)$  a conditional probability distribution over the disturbance space.

We assume that we have a training set  $\mathcal{TS} = \{(i^l, o^l)\}_{l=1}^{\#\mathcal{TS}}$  such that  $o^l$  has been generated from  $i^l$  by the following mechanism: draw  $\xi \in \Xi$  according to  $P_\xi(\cdot|i^l)$  and then set  $o^l = g(i^l, \xi)$ .

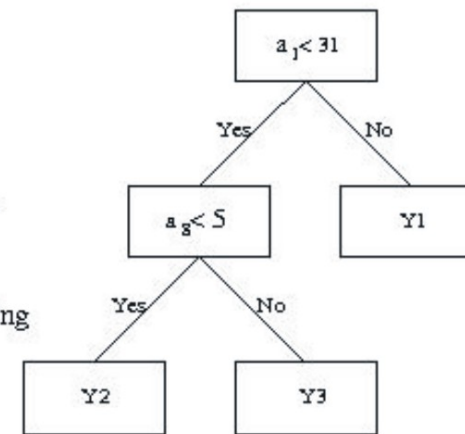
From the sole knowledge of  $\mathcal{TS}$ , supervised learning aims at finding a function  $\hat{g} : I \rightarrow O$  which is a 'good approximation' of the function

$$\bar{g}(i) = E_{\xi \sim P_\xi(\cdot)} [g(i, \xi)]$$

- Typical supervised learning methods are: **kernel-based methods**, **(deep) neural networks**, **tree-based methods**.

$a_1$	$a_2$	$a_3$	$a_4$	$a_5$	$a_6$	$a_7$	$a_8$	Y
60	19	18	17	0	1	1	1	Y1
60	3	22	23	1	29	11	23	Y1
75	9	2	1	3	77	46	3	Y1
2	10	10	2	234	0	0	0	Y2
3	7	9	18	5	0	0	0	Y2
2	14	5	10	8	10	8	10	Y3
65	3	20	21	2	0	1	1	?

⇒  
Batch-mode  
Supervised Learning



- Supervised learning **highly successful**: state-of-the art SL algorithms have been successfully applied to problems where the input state was composed thousands of components.

## The fitted $Q$ iteration algorithm

- **Fitted  $Q$  iteration** computes from  $h_t$  the functions  $\hat{Q}_1, \hat{Q}_2, \dots, \hat{Q}_N$ , approximations of  $Q_1, Q_2, \dots, Q_N$ . At step  $N > 1$ , the algorithm uses the function  $\hat{Q}_{N-1}$  together with  $h_t$  to compute a new training set from which a SL algorithm outputs  $\hat{Q}_N$ . More precisely, this iterative algorithm works as follows:

**First iteration:** the algorithm determines a model  $\hat{Q}_1$  of

$Q_1(x, u) = \underset{w \sim P_w(\cdot|x, u)}{E} [r(x, u, w)]$  by running a SL algorithms on the training set:

$$\mathcal{TS} = \{((x_k, u_k), r_k)\}_{k=0}^{t-1} \quad (16)$$

**Motivation:** One can assimilate  $X \times U$  to  $I$ ,  $\mathbb{R}$  to  $O$ ,  $W$  to  $\Xi$ ,  $P_w(\cdot|x, u)$  to  $P_\xi(\cdot|x, u)$ ,  $r(x, u, w)$  to  $g(i, \xi)$  and  $Q_1(x, u)$  to  $\bar{g}$ . From there, we can observe that a SL algorithm applied to the training set described by equation (16) will produce a model of  $Q_1$ .

**Iteration  $N > 1$ :** the algorithm outputs a model  $\hat{Q}_N$  of  $Q_N(x, u) = E_{w \sim P_w(\cdot|x, u)} [r(x, u, w) + \gamma \max_{u' \in U} Q_{N-1}(f(x, u, w), u')]$  by running a SL algorithms on the training set:

$$\mathcal{TS} = \{((x_k, u_k), r_k + \gamma \max_{u' \in U} \hat{Q}_{N-1}(x_{k+1}, u'))\}_{k=0}^{t-1}$$

**Motivation:** One can reasonably suppose that  $\hat{Q}_{N-1}$  is a sufficiently good approximation of  $Q_{N-1}$  to be consider to be equal to this latter function. Assimilate  $X \times U$  to  $I$ ,  $\mathbb{R}$  to  $O$ ,  $W$  to  $\Xi$ ,  $P_w(\cdot|x, u)$  to  $P_\xi(\cdot|x, u)$ ,  $r(x, u, w)$  to  $g(i, \xi)$  and  $Q_N(x, u)$  to  $\bar{g}$ . From there, we observe that a SL algorithm applied to the training set described by equation (17) will produce a model of  $Q_N$ .

- The algorithm stops when  $N$  is 'large enough' and  $\hat{\mu}_N^*(x) \in \arg \max_{u \in U} \hat{Q}_N(x, u)$  is taken as approximation of  $\mu^*(x)$ .

## The fitted $Q$ iteration algorithm: some remarks

- Performances of the algorithm depends on the supervised learning (SL) method chosen.
- Excellent performances have been observed when combined with supervised learning methods based on **ensemble of regression trees**.
- Fitted  $Q$  iteration algorithm can be used with any set of **one-step system transitions**  $(x_t, u_t, r_t, x_{t+1})$  where each one-step system transition gives information about: a state, the action taken while being in this state, the reward signal observed and the next state reached.
- **Consistency**, that is convergence towards an optimal solution when the number of one-step system transitions tends to infinity, **can be ensured** under appropriate assumptions on the SL method, the sampling process, the system dynamics and the reward function.

## Computation of $\hat{\mu}^*$ : from an inference problem to a problem of computational complexity

- When having at one's disposal only a few one-step system transitions, the main problem is a problem of inference.
- Computational complexity of the fitted  $Q$  iteration algorithm grows with the number  $M$  of one-step system transitions  $(x_k, u_k, r_k, x_{k+1})$  (e.g., it grows as  $M \log M$  when coupled with tree-based methods).
- Above a certain number of one-step system transitions, a problem of computational complexity appears.
- Should we rely on algorithms having less inference capabilities than the 'fitted  $Q$  iteration algorithm' but which are also less computationally demanding to mitigate this problem of computational complexity  $\Rightarrow$  Open research question.

- There is a serious problem plaguing every reinforcement learning algorithm known as the **curse of dimensionality**\*: whatever the mechanism behind the generation of the trajectories and without any restrictive assumptions on  $f(x, u, w)$ ,  $r(x, u, w)$ ,  $X$  and  $U$ , **the number of computer operations required to determine (close-to-) optimal policies tends to grow exponentially with the dimensionality of  $X \times U$ .**
- This exponential growth makes these techniques rapidly computationally impractical when the size of the state-action space increases.
- Many researchers in **reinforcement learning/dynamic programming/optimal control theory** focus their effort on designing algorithms able to **break this curse of dimensionality**. Can deep neural networks give a hope?

\*A term introduced by Richard Bellman (the founder of the DP theory) in the fifties.



## Q-learning with parametric function approximators

Let us extend the Q-learning algorithm to the case where a parametric Q-function of the form  $\tilde{Q}(x, u, a)$  is used:

1. Equation (15) provides us with a desired update for  $\tilde{Q}(x_t, u_t, a)$ , here:  $\delta(x_t, u_t) = r_t + \gamma \max_{u \in U} \tilde{Q}(x_{t+1}, u, a) - \tilde{Q}(x_t, u_t, a)$ , after observing  $(x_t, u_t, r_t, x_{t+1})$ .

2. It follows the following change in parameters:

$$a \leftarrow a + \alpha \delta(x_t, u_t) \frac{\partial \tilde{Q}(x_t, u_t, a)}{\partial a}. \quad (17)$$

Appendix : Algorithmic models for computing the fixed point of a contraction mapping and their application to reinforcement learning.

## Contraction mapping

Let  $B(E)$  be the set of all **bounded real-valued functions** defined on an arbitrary set  $E$ . With every function  $R : E \rightarrow \mathbb{R}$  that belongs to  $B(E)$ , we associate the scalar :

$$\|R\|_{\infty} = \sup_{e \in E} |R(e)|. \quad (18)$$

A **mapping**  $G : B(E) \rightarrow B(E)$  is said to be a **contraction mapping** if there exists a scalar  $\rho < 1$  such that :

$$\|GR - GR'\|_{\infty} \leq \rho \|R - R'\|_{\infty} \quad \forall R, R' \in B(E). \quad (19)$$

## Fixed point

$R^* \in B(E)$  is said to be a *fixed point* of a mapping  $G : B(E) \rightarrow B(E)$  if :

$$GR^* = R^*. \quad (20)$$

If  $G : B(E) \rightarrow B(E)$  is a *contraction mapping* then there exists a unique fixed point of  $G$ . Furthermore if  $R \in B(E)$ , then

$$\lim_{k \rightarrow \infty} \|G^k R - R^*\|_{\infty} = 0. \quad (21)$$

From now on, we assume that:

1.  $E$  is finite and composed of  $n$  elements
2.  $G : B(E) \rightarrow B(E)$  is a contraction mapping whose fixed point is denoted by  $R^*$
3.  $R \in B(E)$ .

## Algorithmic models for computing a fixed point

*All elements of  $R$  are refreshed:* Suppose we have the algorithm that updates at stage  $k$  ( $k \geq 0$ )  $R$  as follows :

$$R \leftarrow GR. \quad (22)$$

The value of  $R$  computed by this algorithm converges to the fixed point  $R^*$  of  $G$ . This is an immediate consequence of equation (21).

*One element of  $R$  is refreshed:* Suppose we have the algorithm that selects at each stage  $k$  ( $k \geq 0$ ) an element  $e \in E$  and updates  $R(e)$  as follows :

$$R(e) \leftarrow (GR)(e) \quad (23)$$

leaving the other components of  $R$  unchanged. If each element  $e$  of  $E$  is selected an infinite number of times then the value of  $R$  computed by this algorithm converges to the fixed point  $R^*$ .

*One element of  $R$  is refreshed and noise introduction:* Let  $\eta \in \mathbb{R}$  be a **noise factor** and  $\alpha \in \mathbb{R}$ . Suppose we have the algorithm that selects at stage  $k$  ( $k \geq 0$ ) an element  $e \in E$  and updates  $R(e)$  according to :

$$R(e) \leftarrow (1 - \alpha)R(e) + \alpha((GR)(e) + \eta) \quad (24)$$

leaving the other components of  $R$  unchanged.

We denote by  $e_k$  the element of  $E$  selected at stage  $k$ , by  $\eta_k$  the noise value at stage  $k$  and by  $R_k$  the value of  $R$  at stage  $k$  and by  $\alpha_k$  the value of  $\alpha$  at stage  $k$ . In order to ease further notations we set  $\alpha_k(e) = \alpha_k$  if  $e = e_k$  and  $\alpha_k(e) = 0$  otherwise.

With this notation equation (24) can be rewritten equivalently as follows :

$$R_{k+1}(e_k) = (1 - \alpha_k)R_k(e_k) + \alpha_k((GR_k)(e_k) + \eta_k). \quad (25)$$

We define the **history**  $\mathcal{F}_k$  of the algorithm at stage  $k$  as being :

$$\mathcal{F}_k = \{R_0, \dots, R_k, e_0, \dots, e_k, \alpha_0, \dots, \alpha_k, \eta_0, \dots, \eta_{k-1}\}. \quad (26)$$

We assume moreover that the following conditions are satisfied:

1. For every  $k$ , we have

$$E[\eta_k | \mathcal{F}_k] = 0. \quad (27)$$

2. There exist two constants  $A$  and  $B$  such that  $\forall k$

$$E[\eta_k^2 | \mathcal{F}_k] \leq A + B \|R_k\|_\infty^2. \quad (28)$$

3. The  $\alpha_k(e)$  are nonnegative and satisfy

$$\sum_{k=0}^{\infty} \alpha_k(e) = \infty, \quad \sum_{k=0}^{\infty} \alpha_k^2(e) < \infty. \quad (29)$$

Then the algorithm **converges with probability 1** to  $R^*$ .

## The $Q$ -function as a fixed point of a contraction mapping

We define the **mapping**  $H: B(X \times U) \rightarrow B(X \times U)$  such that

$$(HK)(x, u) = \underset{w \sim P_w(\cdot|x, u)}{E} [r(x, u, w) + \gamma \max_{u' \in U} K(f(x, u, w), u')] \quad (30)$$

$\forall (x, u) \in X \times U$ .

- The recurrence equation (4) for computing the  $Q_N$ -functions can be rewritten  $Q_N = HQ_{N-1} \forall N > 1$ , with  $Q_0(x, u) \equiv 0$ .
- We prove afterwards that  **$H$  is a contraction mapping**. As immediate consequence, we have, by virtue of the properties algorithmic model (22), that the sequence of  $Q_N$ -functions converges to the unique solution of the Bellman equation (6) which can be rewritten:  $Q = HQ$ . Afterwards, we proof, by using the properties of the algorithmic model (25), the convergence of the  $Q$ -learning algorithm.



## $H$ is a contraction mapping

This  $H$  mapping is a **contraction mapping**. Indeed, we have for any functions  $K, \bar{K} \in B(X \times U)$  :\*  
\*

$$\begin{aligned} \|HK - H\bar{K}\|_\infty &= \gamma \max_{(x,u) \in X \times U} \left| \mathbb{E}_{w \sim P_w(\cdot|x,u)} [\max_{u' \in U} K(f(x,u,w), u') - \right. \\ &\quad \left. \max_{u' \in U} \bar{K}(f(x,u,w), u')] \right| \\ &\leq \gamma \max_{(x,u) \in X \times U} \left| \mathbb{E}_{w \sim P_w(\cdot|x,u)} [\max_{u' \in U} |K(f(x,u,w), u') - \right. \\ &\quad \left. \bar{K}(f(x,u,w), u')|] \right| \\ &\leq \gamma \max_{x \in X} \max_{u \in U} |K(x,u) - \bar{K}(x,u)| \\ &= \gamma \|K - \bar{K}\|_\infty \end{aligned}$$

\*We do as additional assumption here that the rewards are strictly positive.

## Q-learning convergence proof

The Q-learning algorithm updates  $Q$  at stage  $k$  in the following way\*

$$Q_{k+1}(x_k, u_k) = (1 - \alpha_k)Q_k(x_k, u_k) + \alpha_k(r(x_k, u_k, w_k) + \quad (31)$$

$$\gamma \max_{u \in U} Q_k(f(x_k, u_k, w_k), u)), \quad (32)$$

$Q_k$  representing the estimate of the Q-function at stage  $k$ .  $w_k$  is drawn independently according to  $P_w(\cdot | x_k, u_k)$ .

By using the  $H$  mapping definition (equation (30)), equation (32) can be rewritten as follows :

$$Q_{k+1}(x_k, u_k) = (1 - \alpha_k)Q_k(x_k, u_k) + \alpha_k((HQ_k)(x_k, u_k) + \eta_k) \quad (33)$$

\*The element  $(x_k, u_k, r_k, x_{k+1})$  used to refresh the Q-function at iteration  $k$  of the Q-learning algorithm is “replaced” here by  $(x_k, u_k, r(x_k, u_k, w_k), f(x_k, u_k, w_k))$ .

with

$$\begin{aligned}\eta_k &= r(x_k, u_k, w_k) + \gamma \max_{u \in U} Q_k(f(x_k, u_k, w_k), u) - (HQ_k)(x_k, u_k) \\ &= r(x_k, u_k, w_k) + \gamma \max_{u \in U} Q_k(f(x_k, u_k, w_k), u) - \\ &\quad E_{w \sim P_w(\cdot|x, u)} [r(x_k, u_k, w) + \gamma \max_{u \in U} Q_k(f(x_k, u_k, w), u)]\end{aligned}$$

which has exactly the same form as equation (25) ( $Q_k$  corresponding to  $R_k$ ,  $H$  to  $G$ ,  $(x_k, u_k)$  to  $e_k$  and  $X \times U$  to  $E$ ).

We know that  $H$  is a contraction mapping. If the  $\alpha_k(x_k, u_k)$  terms satisfy expression (29), we still have to verify that  $\eta_k$  satisfies expressions (27) and (28), where

$$\mathcal{F}_k = \{Q_0, \dots, Q_k, (x_0, u_0), \dots, (x_k, u_k), \alpha_0, \dots, \alpha_k, \eta_0, \dots, \eta_{k-1}\}, \quad (34)$$

in order to ensure the convergence of the Q-learning algorithm.

We have :

$$\begin{aligned} E[\eta_k | \mathcal{F}_k] &= E_{w_k \sim P_w(\cdot | x_k, u_k)} [r(x_k, u_k, w_k) + \gamma \max_{u \in U} Q_k(f(x_k, u_k, w_k), u) - \\ &\quad E_{w \sim P_w(\cdot | x_k, u_k)} [r(x_k, u_k, w) + \gamma \max_{u \in U} Q_k(f(x_k, u_k, w), u)] | \mathcal{F}_k] \\ &= 0 \end{aligned}$$

and expression (27) is indeed satisfied.

In order to prove that expression (28) is satisfied, one can first note that :

$$|\eta_k| \leq 2B_r + 2\gamma \max_{(x,u) \in X \times U} Q_k(x, u) \quad (35)$$

where  $B_r$  is the bound on the rewards. Therefore we have :

$$\eta_k^2 \leq 4B_r^2 + 4\gamma^2 \left( \max_{(x,u) \in X \times U} Q_k(x, u) \right)^2 + 8B_r\gamma \max_{(x,u) \in X \times U} Q_k(x, u) \quad (36)$$

By noting that

$$8B_r\gamma \max_{(x,u) \in X \times U} Q_k(x, u) < 8B_r\gamma + 8B_r\gamma \left( \max_{(x,u) \in X \times U} Q_k(x, u) \right)^2 \quad (37)$$

and by choosing  $A = 8B_r\gamma + 4B_r^2$  and  $B = 8B_r\gamma + 4\gamma^2$  we can write

$$\eta_k^2 \leq A + B \|Q_k\|_\infty^2 \quad (38)$$

and expression (28) is satisfied. **QED**

## Additional readings

Lucian Busoniu, Robert Babuska, Bart De Schutter, and Damien Ernst. Reinforcement learning and dynamic programming using function approximators. CRC Press, April 2010. (Available at: <https://orbi.ulg.ac.be/bitstream/2268/27963/1/book-FA-RL-DP.pdf>)

Richard Sutton and Andrew Barto. Reinforcement Learning: An Introduction. Second edition, in progress. Available at: <http://ufal.mff.cuni.cz/~straka/courses/npfl114/2016/sutton-bookdraft2016sep.pdf>)

Dimitri P. Bertsekas. Dynamic Programming and Optimal Control. Volume I (last edition: 2017) and Volume II (last edition: 2012). Material available at <http://ufal.mff.cuni.cz/~straka/courses/npfl114/2016/sutton-bookdraft2016sep.pdf>

Csaba Szepesvari. Algorithms for Reinforcement Learning. Synthesis Lectures on Artificial Intelligence and Machine Learning. June 2009. Available at <http://jmlr.csail.mit.edu/papers/v6/ernst05a.html>)