

Neuro & AI: Methods and Tools for Neuroergonomics

Introduction to Reinforcement Learning

Giorgio Angelotti^{1,2} **Nicolas Drougard¹**

¹ISAE-SUPAERO DCAS, Toulouse, FRANCE

²ANITI, Toulouse, FRANCE

`nicolas.drougard@isae-sup aero.fr`

- 1 Introduction
- 2 Markov Decision Processes
- 3 Reinforcement Learning

- ▶ [HMS20] G. Haine, D. Matignon, and M. Salaün, *Mathématiques déterministes*, Tronc Commun Scientifique 1A, Formation Ingénieur ISAE-SUPAERO, 2020.
- ▶ [MKS⁺15] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A. Rusu, Joel Veness, Marc G. Bellemare, Alex Graves, Martin Riedmiller, Andreas K. Fidjeland, Georg Ostrovski, Stig Petersen, Charles Beattie, Amir Sadik, Ioannis Antonoglou, Helen King, Dhharshan Kumaran, Daan Wierstra, Shane Legg, and Demis Hassabis, *Human-level control through deep reinforcement learning*, Nature **518** (2015), no. 7540, 529–533, download link [here](#).
- ▶ [Mun13] Rémi Munos, *Lectures: Apprentissage par renforcement*, 2013, download links [here](#).
- ▶ [SB18] Richard S Sutton and Andrew G Barto, *Reinforcement learning: An introduction*, MIT press, 2018, download link [here](#).
- ▶ [Sil15] David Silver, *Lectures on reinforcement learning*, 2015, download links [here](#).
- ▶ [Zan05] Elena Zanini, *Markov decision processes*, 2005, download link [here](#).

- **The** book on Reinforcement Learning (RL) [SB18]
- Lectures on RL by David Silver [Sil15]
- Short presentation of Markov Decision Processes [Zan05]
- French lectures on RL by R.Munos [Mun13]

- 1 Introduction
- 2 Markov Decision Processes
- 3 Reinforcement Learning

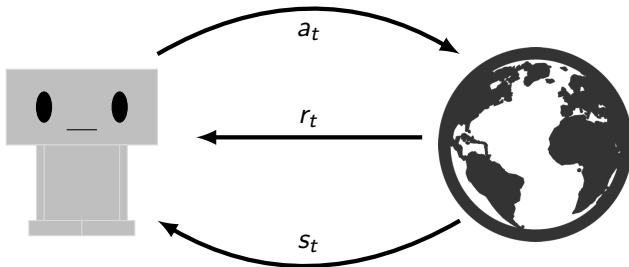
- RL = Reinforcement Learning
- Learn a behaviour from interaction (with a simulator).
- No supervision (label $y \in \mathcal{Y}$), only a reward signal ($R_t \in \mathbb{R}$).
- Growing field, benefits from Deep Learning.
- Classical RL operates on finite spaces, Deep RL handles continuous ones.
- Videos: Breakout, Locomotion behavior, Alphastar.
- ISAE-SUPAERO RL Initiative SuReLI.

Classical Reinforcement learning is based on the **reward hypothesis**.

Reward Hypothesis

A goal can be described by the maximisation of expected cumulative reward $\mathbb{E} [\sum_{t \geq 0} R_t]$.

Do you agree with this statement?



Reinforcement Learning?

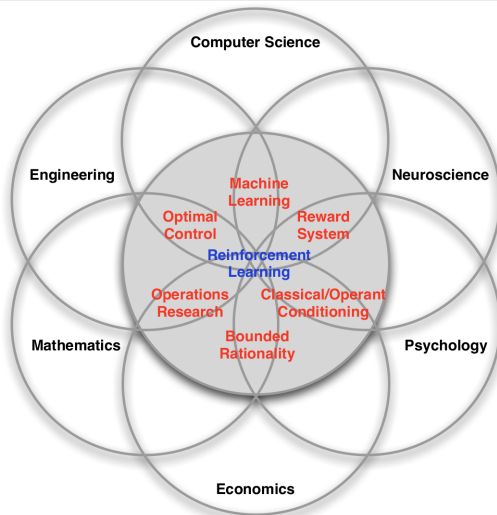


Illustration from the first lecture of David Silver [Sil15].

- 1 Introduction
- 2 Markov Decision Processes
- 3 Reinforcement Learning

The origin of Markov Processes

- Random model set up by Andreï Markov (1856-1922).



- In particular, he studied the letter sequence process in novels.
- Each letter depends primarily on the previous one.

Definition – Markov Chain (or Process)

- \mathcal{S} countable set called *set of states* $s \in \mathcal{S}$.
- $(S_t)_{t \in \mathbb{N}}$ sequence of random variables of \mathcal{S} .

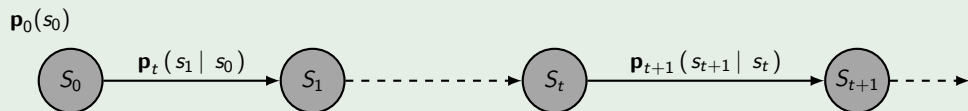
The sequence $(S_t)_{t \in \mathbb{N}}$ is a **Markov Chain** if $\forall t \geq 0, \forall (s_0, s_1, \dots, s_{t+1}) \in \mathcal{S}^{t+2}$

$$\mathbb{P}(S_{t+1} = s_{t+1} \mid S_0 = s_0, S_1 = s_1, \dots, S_t = s_t) = \mathbb{P}(S_{t+1} = s_{t+1} \mid S_t = s_t).$$

Other equivalent definitions

- “the future is independent of the past given the present”
- $\forall t \geq 1, S_{t+1} \perp\!\!\!\perp \{S_i\}_{i=0}^{t-1} \mid S_t$.
note: “ $X \perp\!\!\!\perp Y \mid Z$ ” means “ X is independent from Y given Z ”.
- $S_0 \sim P_0$ and $\forall t > 0, S_t \sim P_{s,t}$ if $S_{t-1} = s$,
where p_0 and $\{P_{s,t}\}_{(s,t) \in \mathcal{S} \times \mathbb{N}}$ are probability distributions on \mathcal{S} .
note: “ $X \sim P$ ” means “ X has the probability distribution P ”.

Graphical representation of a Markov Process



- Graphical representation called **Bayesian Network** (BN).
- In BNs, for each set of nodes N ,

$$N \perp\!\!\!\perp \text{non-descendant}(N) \mid \text{parents}(N).$$

- Here, $\text{parents}(S_t) = \{S_{t-1}\}$, $\text{non-descendant}(S_t) = \{S_0, S_1, \dots, S_{t-2}\}$.
- For a better readability, $\mathbb{P}(S_{t+1} = s' \mid S_t = s)$ is denoted by $\mathbf{p}_t(s' \mid s)$: $P_{s,t} = \mathbf{p}(\cdot \mid s)$.

Definition – Transition matrix

Let's assume that $\#\mathcal{S} = n$ and let's assign a number to each state $\mathcal{S} = \{s^{(1)}, s^{(2)}, \dots, s^{(n)}\}$.

The transition matrix is $\forall (i, j) \in \{1, \dots, n\}^2, \forall t \geq 0, T_{ij}^{(t)} = \mathbf{p}_t(s^{(j)} | s^{(i)})$, i.e.

$$T^{(t)} = \begin{pmatrix} \mathbf{p}_t(s^{(1)} | s^{(1)}) & \dots & \mathbf{p}_t(s^{(n)} | s^{(1)}) \\ \mathbf{p}_t(s^{(1)} | s^{(2)}) & \dots & \mathbf{p}_t(s^{(n)} | s^{(2)}) \\ \vdots & & \vdots \\ \mathbf{p}_t(s^{(1)} | s^{(n)}) & \dots & \mathbf{p}_t(s^{(n)} | s^{(n)}) \end{pmatrix} = \begin{pmatrix} P_{s^{(1)}, t} \\ P_{s^{(2)}, t} \\ \vdots \\ P_{s^{(n)}, t} \end{pmatrix}.$$

Distribution of S_t

Consider $p_t \in \mathbb{R}^n$ such that $\forall i \in \{1, \dots, n\}, (p_t)_i = \mathbb{P}(S_t = s^{(i)})$. Thus, $\forall t \geq 0$,

$$\begin{aligned} (p_t)_i &= \sum_{j=1}^n \mathbb{P}(S_t = s^{(i)}, S_{t-1} = s^{(j)}) = \sum_{j=1}^n \mathbf{p}_t(s^{(i)} | s^{(j)}) \mathbb{P}(S_{t-1} = s^{(j)}) \\ &= (T^{(t)} p_{t-1})_i = (T^{(t)} T^{(t-1)} p_{t-2})_i = \dots = (T^{(t)} \dots T^{(1)} p_0)_i. \end{aligned}$$

Property of Markov Processes (MPs)

Let $(S_t)_{t \in \mathbb{N}}$ be a Markov Process on space \mathcal{S} .

$\forall t \in \mathbb{N}, \forall r : \mathcal{S} \rightarrow \mathbb{R}$ bounded, $\forall (s_0, \dots, s_t) \in \mathcal{S}^{t+1}$,

$$\begin{aligned}\mathbb{E}[r(S_{t+1}) \mid S_0 = s_0, \dots, S_t = s_t] &= \sum_{s' \in \mathcal{S}} r(s') \cdot \mathbb{P}(S_{t+1} = s' \mid S_t = s_t) \\ &= \mathbb{E}[r(S_{t+1}) \mid S_t = s_t].\end{aligned}$$

Proof

$$\begin{aligned}\mathbb{E}[r(S_{t+1}) \mid S_0 = s_0, \dots, S_t = s_t] &= \mathbb{E}\left[\sum_{s' \in \mathcal{S}} r(s') \cdot \mathbb{1}_{\{S_{t+1}=s'\}} \mid S_0 = s_0, \dots, S_t = s_t\right] \\ \left(\mathbb{1}_{\{e\}} \stackrel{\text{def}}{=} \begin{cases} 1 & \text{if } e \text{ true,} \\ 0 & \text{otherwise.} \end{cases}\right) &\quad \stackrel{\mathbb{E} \text{ linear}}{=} \sum_{s' \in \mathcal{S}} r(s') \cdot \mathbb{E}\left[\mathbb{1}_{\{S_{t+1}=s'\}} \mid S_0 = s_0, \dots, S_t = s_t\right] \\ &\stackrel{\mathbb{P} \text{ def}}{=} \sum_{s' \in \mathcal{S}} r(s') \cdot \mathbb{P}(S_{t+1} = s' \mid S_0 = s_0, \dots, S_t = s_t) \\ &\stackrel{MP}{=} \sum_{s' \in \mathcal{S}} r(s') \cdot \mathbb{P}(S_{t+1} = s' \mid S_t = s_t) = \mathbb{E}[r(S_{t+1}) \mid S_t = s_t].\end{aligned}$$

Decision Process (DP)

- A DP models a *system* depending on *actions* over *time*.
- Possible *states of the system* are $s \in \mathcal{S}$.
- \mathbb{N} models the *time*, or the *stages of the process*.
- Possible *actions* are denoted by “ a ” chosen from a finite set \mathcal{A} .
- The *agent* = the entity responsible for decision making
i.e. choosing the successive actions given the current information about the system.
- Example: *DP for a robot*
 - System state $s \in \mathcal{S}$ features of both the robot and its environment needed to describe the robotic mission, as the robot location or the target location.
 - Actions executable by the robot $a \in \mathcal{A}$.
 - Agent = decision making module of the robot selecting the robot's actions.

Definition – Markov Decision Process (MDP)

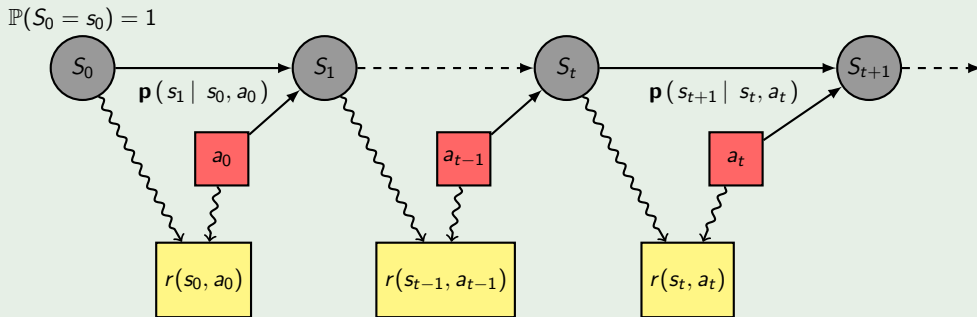
An MDP is a 4-uple $\langle \mathcal{S}, \mathcal{A}, T, r \rangle$, where

- \mathcal{S} is the finite set of **system states** $s \in \mathcal{S}$,
- \mathcal{A} is the finite set of available **actions** $a \in \mathcal{A}$,
- T is the **transition function**: $T(s, a, s') \stackrel{\text{def}}{=} \mathbf{p}(s' \mid s, a)$, $\forall (s, s', a) \in \mathcal{S}^2 \times \mathcal{A}$,
- $r : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{R}$ is the **reward function**.

Remark

If we define a **plan** $(a_t)_{t \in \mathbb{N}}$, i.e. a sequence of actions $a_t \in \mathcal{A}$, S_t with transition matrices $T_{ij}^{(t)} = \mathbf{p}(s^{(j)} \mid s^{(i)}, a_t)$ is a Markov Chain.

Graphical representation of a Markov Decision Process



- The gray circles illustrate the Markov process of system states,
- the red squares are decisions, selected actions over time,
- and the yellow squares are collected rewards.

Definition – Strategy

A **strategy** (or policy) is a sequence of functions $\pi = (\pi_t)_{t \in \mathbb{N}}$ such that $\pi_t : \mathcal{S}^{t+1} \rightarrow \mathcal{A}$. Given a strategy π , \forall state **history** $h_t = (s_0, \dots, s_t) \in \mathcal{S}^{t+1}$, $\forall t \in \mathbb{N}$, the selected action is

$$a_t = \pi_t(s_0, s_1, \dots, s_t) = \pi_t(h_t).$$

Solving an MDP

- Given a horizon $H \in \mathbb{N}$, the criterion of an MDP is the **Value Function**:

$$V_H^\pi(s) = \mathbb{E} \left[\sum_{t=0}^H r(S_t, \pi_t(h_t)) \mid S_0 = s \right].$$

- A solution of an MDP is a strategy $(\pi_t^*)_{t \in \mathbb{N}}$ maximizing the value function:

$$\pi^* \in \arg \max_{\pi} V^\pi(s).$$

Dynamic Programming

To solve the finite horizon problem ($H < +\infty$), let's start from the end:

- Decision step H : what is the optimal choice π_H^* ?

- the agent knows $\{S_0 = s_0, \dots, S_H = s_H\}$,

- whatever the previous decisions $(\pi_t)_{t=0}^{H-1}$, π_H^* should maximize

$$a \mapsto \sum_{t=0}^{H-1} r(s_t, \pi_t(s_t)) + r(s_H, a) = cst + r(s_H, a).$$

$$\Rightarrow \pi_H^* \text{ depends only on } s_H: \pi_H^*(s_H) = \arg \max_{a \in \mathcal{A}} (r(s_H, a)) = \arg \max_{a \in \mathcal{A}} (v_H(s_H, a)),$$

with $v_H(s, a) = r(s, a)$.

The optimal last decision step π_H^* is known.

Let's define $v_H^*(s) = \max_{a \in \mathcal{A}} v_H(s, a) = \max_{a \in \mathcal{A}} r(s, a)$,

and consider the previous decision step $H - 1$.

Dynamic Programming

■ Decision step $H - 1$: what is the optimal choice π_{H-1}^* ?

- the agent knows $\{S_0 = s_0, \dots, S_{H-1} = s_{H-1}\}$, and the optimal last decision π_H^* .
- whatever the previous decisions $(\pi_t)_{t=0}^{H-2}$, π_{H-1}^* should maximize

$$\begin{aligned}
 a \mapsto & \sum_{t=0}^{H-2} r(s_t, \pi_t(s_t)) + r(s_{H-1}, a) + \mathbb{E}[r(S_H, \pi_H^*(S_H)) \mid S_{H-1} = s_{H-1}] \\
 & = cst + r(s_{H-1}, a) + \mathbb{E}[r(S_H, \pi_H^*(S_H)) \mid S_{H-1} = s_{H-1}] \\
 & = cst + r(s_{H-1}, a) + \sum_{s' \in \mathcal{S}} T(s_{H-1}, a, s') r(s', \pi_H^*(s')) \\
 & = cst + r(s_{H-1}, a) + \sum_{s' \in \mathcal{S}} T(s_{H-1}, a, s') v_H^*(s').
 \end{aligned}$$

\Rightarrow given π_H^* already computed, π_{H-1}^* depends only on s_{H-1} :

$$\pi_{H-1}^*(s_{H-1}) = \arg \max_{a \in \mathcal{A}} \left(r(s_{H-1}, a) + \sum_{s' \in \mathcal{S}} T(s_{H-1}, a, s') v_H^*(s') \right) = \arg \max_{a \in \mathcal{A}} (v_{H-1}(s_{H-1}, a)),$$

where $v_H^*(s) = \max_{a \in \mathcal{A}} v_H(s, a)$ and $v_{H-1}(s, a) = r(s, a) + \sum_{s' \in \mathcal{S}} T(s, a, s') v_H^*(s')$.

The optimal two last decisions are known: π_H^* and π_{H-1}^* .

We define $v_{H-1}^*(s) = \max_{a \in \mathcal{A}} v_{H-1}(s, a)$, and consider the previous step $H - n$, with $2 \leq n \leq H$.

Dynamic Programming

■ Decision step $H - n$: what is the optimal choice π_{H-n}^* ?

- the agent knows $\{S_0 = s_0, \dots, S_{H-n} = s_{H-n}\}$, and optimal decisions from step $H - n + 1$.
- whatever the previous decisions $(\pi_t)_{t=0}^{H-n-1}$, π_{H-n}^* should maximize

$$\begin{aligned} a \mapsto & \sum_{t=0}^{H-n-1} r(s_t, \pi_t(s_t)) + r(s_{H-n}, a) + \mathbb{E} \left[\sum_{t=H-n+1}^H r(s_t, \pi_t^*(s_t)) \mid S_{H-n} = s_{H-n} \right] \\ & = cst + r(s_{H-n}, a) + \mathbb{E} \left[\sum_{t=H-n+1}^H r(s_t, \pi_t^*(s_t)) \mid S_{H-n} = s_{H-n} \right] \\ & = cst + r(s_{H-n}, a) + \sum_{s' \in \mathcal{S}} T(s_{H-n}, a, s') v_{H-n+1}^*(s'). \end{aligned}$$

$\Rightarrow \pi_{H-n}^*$ depends only on s_{H-n} :

$$\pi_{H-n}^*(s_{H-n}) = \arg \max_{a \in \mathcal{A}} \left(r(s_{H-n}, a) + \sum_{s' \in \mathcal{S}} T(s_{H-n}, a, s') v_{H-n+1}^*(s') \right),$$

and we produce the value needed for the previous decision step optimization:

$$v_{H-n-1}^*(s) = \max_{a \in \mathcal{A}} \left(r(s_{H-n}, a) + \sum_{s' \in \mathcal{S}} T(s_{H-n}, a, s') v_{H-n+1}^*(s') \right).$$

Dynamic Programming

We have just shown by induction that

- π_t^* only depends on s_t
- we can compute an optimal strategy as follows:

$$\pi_H^*(s) = \arg \max_{a \in \mathcal{A}} (r(s, a)),$$

$$v_H^*(s) = \max_{a \in \mathcal{A}} (r(s, a)),$$

and for t from $H - 1$ to 0 :

$$\pi_t^*(s) = \arg \max_{a \in \mathcal{A}} \left(r(s, a) + \sum_{s' \in \mathcal{S}} T(s, a, s') v_{t+1}^*(s') \right),$$

$$v_t^*(s) = \max_{a \in \mathcal{A}} \left(r(s, a) + \sum_{s' \in \mathcal{S}} T(s, a, s') v_{t+1}^*(s') \right).$$

Consequences on strategies

A strategy can then be defined as

$$\pi_t : \mathcal{S} \rightarrow \mathcal{A}.$$

Note: in the litterature, it can also return a probability distribution on \mathcal{A} .

Consequence on MDPs

Given a strategy π , $(S_t)_{t \in \mathbb{N}}$ is a Markov Chain with transition matrices

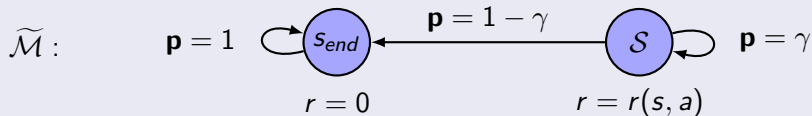
$$T_{ij}^{(t)} = \mathbf{p} \left(s^{(j)} \mid s^{(i)}, \pi_t(s^{(i)}) \right).$$

Introduction to the infinite horizon case ($H = +\infty$)

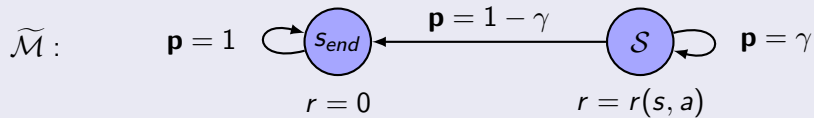
Issue with case $H = +\infty$: $\sum_{t=0}^{+\infty} r(s_t, a_t)$ not always finite.

Overcoming this issue by defining a new MDP $\widetilde{\mathcal{M}}$ with an extra state s_{end} .

- s_{end} absorbing state (or dead end): $\mathbf{p}(s_{end} | s_{end}, a) = 1, \forall a \in \mathcal{A}$.
- $r(s_{end}, a) = 0, \forall a \in \mathcal{A}$.
- $\mathbf{p}(s_{end} | s, a) = 1 - \gamma \in]0, 1[, \forall (s, a) \in \mathcal{S} \times \mathcal{A}$.



Introduction to the infinite horizon case ($H = +\infty$)



From the initial MDP $\mathcal{M} = \langle \mathcal{S}, \mathcal{A}, T, r \rangle$, a new MDP $\tilde{\mathcal{M}} = \langle \tilde{\mathcal{S}}, \mathcal{A}, \tilde{T}, \tilde{r} \rangle$ is defined, with:

- a system space $\tilde{\mathcal{S}} = \mathcal{S} \cup \{s_{end}\}$.

- a new transition function $\tilde{T}(s, a, s') = \begin{cases} \gamma T(s, a, s') & \text{if } (s, s') \in \mathcal{S}^2 \\ (1 - \gamma) & \text{if } s \in \mathcal{S}, s' = s_{end}, \\ 1 & \text{if } s = s' = s_{end}, \\ 0 & \text{if } s = s_{end}, s' \in \mathcal{S}. \end{cases}$

- a new reward function $\tilde{r}(s, a) = \begin{cases} r(s, a) & \text{if } s \in \mathcal{S}, \\ 0 & \text{if } s = s_{end}. \end{cases}$

Value function of the new MDP $\tilde{\mathcal{M}}$

If $S_0 = s_{end}$, $S_t = s_{end} \forall t \geq 0$: $\tilde{\mathcal{M}}$'s value function is $\tilde{V}_H^\pi(s_{end}) = \sum_{t=0}^H r(s_{end}, \pi_t(h_t)) = 0$.

Consider now $s \in \mathcal{S}$, i.e. $s \neq s_{end}$:

$$\tilde{V}_0^\pi(s) = \mathbb{E}_{S_t \sim \tilde{T}} \left[r(S_0, \pi_0(S_0)) \mid S_0 = s \right] = r(s, \pi_0(s)).$$

$$\begin{aligned} \tilde{V}_1^\pi(s) - \tilde{V}_0^\pi(s) &= \mathbb{E}_{S_t \sim \tilde{T}} \left[r(S_1, \pi_1(S_1)) \mid S_0 = s \right] = \sum_{s' \in \tilde{\mathcal{S}}} \tilde{T}(s, \pi_1(s'), s') r(s', \pi_1(s')) \\ &= (1 - \gamma) r(s_{end}, \pi_1(s_{end})) + \sum_{s' \in \mathcal{S}} \gamma T(s, \pi_1(s'), s') r(s', \pi_1(s')) \\ &= 0 + \gamma \sum_{s' \in \mathcal{S}} T(s, \pi_1(s'), s') r(s', \pi_1(s')) \\ &= \mathbb{E}_{S_t \sim T} \left[\gamma r(S_1, \pi_1(S_1)) \mid S_0 = s \right]. \end{aligned}$$

Value function of the new MDP $\widetilde{\mathcal{M}}$

Using strategy $(\pi_t)_{t \in N}$, and considering $f : \mathcal{S} \rightarrow \mathbb{R}$ with $f(s_{end}) = 0$:

$$\begin{aligned}\mathbb{E}_{S_{t+1} \sim \widetilde{T}} \left[f(S_{t+1}) \mid S_t = s \right] &= \sum_{s' \in \widetilde{\mathcal{S}}} \widetilde{T}(s, \pi_t(s), s') f(s') \\ &= (1 - \gamma) f(s_{end}) + \sum_{s' \in \mathcal{S}} \gamma T(s, \pi_1(s'), s') f(s') \\ &= 0 + \gamma \sum_{s' \in \mathcal{S}} T(s, \pi_1(s'), s') f(s') \\ &= \gamma \mathbb{E}_{S_{t+1} \sim T} \left[f(S_{t+1}) \mid S_t = s \right].\end{aligned}$$

Some properties of the conditional expectation

$X \in \mathcal{X}$, $Y \in \mathcal{Y}$ random variables, $f : \mathcal{X} \rightarrow \mathbb{R}$ a function of X returning real numbers.

- $\mathbb{E}[f(X) \mid Y]$ is a function of Y :

$$\sum_{x \in \mathcal{X}} f(x) \mathbb{P}(X = x \mid Y).$$

- $\mathbb{E}[f(X) \mid X] = f(X).$
- $\mathbb{E}[\mathbb{E}[f(X) \mid Y] \mid Z] = \mathbb{E}[f(X) \mid Z].$

Value function of the new MDP $\tilde{\mathcal{M}}$

$$\begin{aligned}
 \tilde{V}_H^\pi(s) - \tilde{V}_{H-1}^\pi(s) &= \mathbb{E}_{S_t \sim \tilde{T}} \left[r(S_H, \pi_H(S_H)) \mid S_0 = s \right] \\
 &= \mathbb{E}_{S_t \sim \tilde{T}} [f(S_H) \mid S_0 = s] \\
 &= \mathbb{E}_{S_t \sim \tilde{T}} \left[\mathbb{E}_{S_t \sim \tilde{T}} [f(S_H) \mid S_{H-1}] \mid S_0 = s \right] \\
 &= \mathbb{E}_{S_t \sim \tilde{T}} [\gamma \mathbb{E}_{S_t \sim T} [f(S_H) \mid S_{H-1}] \mid S_0 = s] \\
 &= \mathbb{E}_{S_t \sim \tilde{T}} [\gamma g(S_{H-1}) \mid S_0 = s] \\
 &= \mathbb{E}_{S_t \sim \tilde{T}} \left[\gamma \mathbb{E}_{S_t \sim \tilde{T}} [g(S_{H-1}) \mid S_{H-2}] \mid S_0 = s \right] \\
 &= \mathbb{E}_{S_t \sim \tilde{T}} \left[\gamma^2 \mathbb{E}_{S_t \sim T} [g(S_{H-1}) \mid S_{H-2}] \mid S_0 = s \right] \\
 &= \mathbb{E}_{S_t \sim \tilde{T}} \left[\gamma^2 \mathbb{E}_{S_t \sim T} [f(S_H) \mid S_{H-2}] \mid S_0 = s \right] \\
 &= \dots \\
 &= \mathbb{E}_{S_t \sim \tilde{T}} \left[\gamma^H \mathbb{E}_{S_t \sim T} [f(S_H) \mid S_0] \mid S_0 = s \right] = \gamma^H \mathbb{E}_{S_t \sim T} [f(S_H) \mid S_0 = s].
 \end{aligned}$$

Value function of the new MDP $\tilde{\mathcal{M}}$

$$\begin{aligned}\tilde{V}_H^\pi(s) &= \sum_{t=0}^H \tilde{V}_t^\pi(s) - \tilde{V}_{t-1}^\pi(s) = \sum_{t=0}^H \gamma^t \mathbb{E}_{S_t \sim T} \left[r(S_t, \pi_t(S_t)) \mid S_0 = s \right] \\ &= \mathbb{E}_{S_t \sim T} \left[\sum_{t=0}^H \gamma^t r(S_t, \pi_t(S_t)) \mid S_0 = s \right].\end{aligned}$$

Since $\sum_{t \geq 0} \gamma^t = \frac{1}{1-\gamma}$, $\sum_{t=0}^H \gamma^t r(S_t, \pi_t(S_t)) \leq \frac{1}{1-\gamma} \max_{(s,a) \in \mathcal{S} \times \mathcal{A}} r(s, a)$,

and the value function of $\tilde{\mathcal{M}}$ can be computed for $H \rightarrow +\infty$.

Value function for an infinite horizon

$$V^\pi(s) = \mathbb{E} \left[\sum_{t \geq 0} \gamma^t r(S_t, \pi_t(S_t)) \mid S_0 = s \right], \text{ with } 0 < \gamma < 1 \text{ (discount factor).}$$

Bellman Equation – Let $\pi_{t \in \mathbb{N}}$ be a strategy.

$$\begin{aligned}
 V^\pi(s) &:= \mathbb{E} \left[\sum_{t=0}^{+\infty} \gamma^t \cdot r(S_t, \pi_t(S_t)) \mid S_0 = s \right] \\
 &= \mathbb{E} \left[r(S_0, \pi_0(S_0)) + \sum_{t=1}^{+\infty} \gamma^t \cdot r(S_t, \pi_t(S_t)) \mid S_0 = s \right] \\
 &= r(s, \pi_0(s)) + \mathbb{E} \left[\mathbb{E} \left[\sum_{t=1}^{+\infty} \gamma^t \cdot r(S_t, \pi_t(S_t)) \mid S_1 \right] \mid S_0 = s \right] \\
 &= r(s, \pi_0(s)) + \gamma \cdot \sum_{s' \in \mathcal{S}} \mathbf{p}(s' | s, \pi_0(s)) \cdot \mathbb{E} \left[\sum_{t=1}^{+\infty} \gamma^{t-1} \cdot r(S_t, \pi_t(S_t)) \mid S_1 = s' \right] \\
 &= r(s, \pi_0(s)) + \gamma \cdot \sum_{s' \in \mathcal{S}} \mathbf{p}(s' | s, \pi_0(s)) \cdot \mathbb{E} \left[\sum_{t'=0}^{+\infty} \gamma^{t'} \cdot r(S_{t'}^+, \pi_{t'}^+(S_{t'}^+)) \mid S_0^+ = s' \right] \\
 &= r(s, \pi_0(s)) + \gamma \cdot \sum_{s' \in \mathcal{S}} \mathbf{p}(s' | s, \pi_0(s)) \cdot V^{\pi^+}(s'), \text{ with } S_t^+ = S_{t+1}, \pi_t^+(s) = \pi_{t+1}(s).
 \end{aligned}$$

Bellman Equation

For a strategy $(\pi_t)_{t \geq 0}$,

$$V^\pi(s) = r(s, \pi_0(s)) + \gamma \cdot \sum_{s' \in \mathcal{S}} \mathbf{p}(s'|s, \pi_0(s)) \cdot V^{\pi^+}(s'),$$

with $\forall s \in \mathcal{S}, \forall t \geq 0, \pi_t^+(s) = \pi_{t+1}(s)$.

Dynamic Programming Equation

The optimal value function $V^* = V^{\pi^*}$ is the unique solution of the following equation:

$$V^*(s) = \max_{a \in \mathcal{A}} \left(r(s, a) + \gamma \cdot \sum_{s' \in \mathcal{S}} \mathbf{p}(s'|s, a) \cdot V^*(s') \right).$$

There exists an optimal strategy $\pi^* : \mathcal{S} \rightarrow \mathcal{A}$, independent from the time step $t \in \mathbb{N}$, that is $\pi^*(s) = \arg \max_{a \in \mathcal{A}} (r(s, a) + \gamma \cdot \sum_{s' \in \mathcal{S}} \mathbf{p}(s'|s, a) \cdot V^*(s'))$.

Q-value

$$Q^{\pi}(s, a) = r(s, a) + \gamma \sum_{s' \in \mathcal{S}} \mathbf{p}(s' \mid s, a) V^{\pi}(s').$$

Optimal Q-value

$$Q^*(s, a) = r(s, a) + \gamma \sum_{s' \in \mathcal{S}} \mathbf{p}(s' \mid s, a) V^*(s').$$

$$\Rightarrow \pi^*(s) = \arg \max_{a \in \mathcal{A}} Q^*(s, a).$$

Value Iteration Algorithm

Initialize: $Q = r$

Repeat:

- $V(s) = \max_{a \in \mathcal{A}} Q(s, a)$
- $Q(s) = r(s, a) + \gamma \sum_{s' \in \mathcal{S}} \mathbf{p}(s' | s, a) V(s')$

(i.e. DP equations) until convergence

(thanks to the Fixed Point/Contraction Mapping Theorem (Thm 7.1 [HMS20]):
Bellman & DP operators are contracting).

Return:

- $\pi^*(s) = \max_a Q(s, a).$

With matrices

Denote states by $\mathcal{S} = \{s^{(1)}, s^{(2)}, \dots, s^{(n)}\}$, with $n = \#\mathcal{S}$.

Define

- action vector $a \in \mathcal{A}^n$,
- strategy vector $\pi \in \mathcal{A}^n$: $\pi_i = \pi(s^{(i)})$, $\forall i \in \{1, \dots, n\}$,
- reward vector $r^a \in \mathbb{R}^n$: $r_i^a = r(a_i, s^{(i)})$, $\forall i \in \{1, \dots, n\}$,
- transition matrix $T^a \in \mathbb{R}^{n \times n}$: $T_{ij}^a = \mathbf{p}(s^{(j)} \mid s^{(i)}, a_i)$, $\forall (i, j) \in \{1, \dots, n\}^2$,
- value vector $V \in \mathbb{R}^n$: $V_i = V(s^{(i)})$, $\forall i \in \{1, \dots, n\}$,

The Bellman Equation becomes

$$V^\pi = r^\pi + \gamma T^\pi V^\pi \Rightarrow (I_n - \gamma T^\pi) V^\pi = r^\pi.$$

The Dynamic Programming Equation becomes

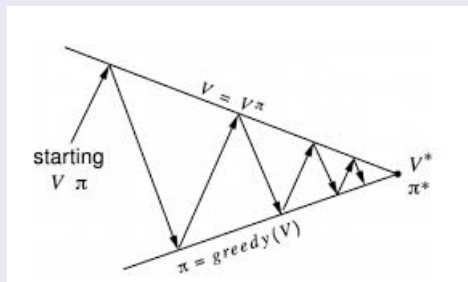
$$V^* = \max_{a \in \mathcal{A}^n} r^a + \gamma T^a V^*.$$

Policy Iteration Algorithm

Initialize $Q = r$, $\pi = a \in \mathcal{A}^n$, and repeat:

- $\pi \in \arg \max_{a \in \mathcal{A}^n} Q$,
- resolution of $(I_n - \gamma T^\pi)V = r^\pi$ (or iterations of Bellman equation)
- $Q = r^\pi + \gamma T^\pi V$.

until convergence (i.e. when π remains the same).



Policy Iteration Algorithm always converges to π^*

Given π , $\forall s \in \mathcal{S}$, $\pi'(s) = \arg \max_{a \in \mathcal{A}} Q^\pi(s, a)$ is a better action: $V^\pi(s) = Q^\pi(s, \pi(s))$

$$\leq \max_{a \in \mathcal{A}} Q^\pi(s, a) = Q^\pi(s, \pi'(s))$$

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

$$= r(s, \pi'(s)) + \gamma \sum_{s' \in \mathcal{S}} T(s, \pi'(s), s') Q^\pi(s', \pi(s'))$$

$$\leq r(s, \pi'(s)) + \gamma \sum_{s' \in \mathcal{S}} T(s, \pi'(s), s') Q^\pi(s', \pi'(s'))$$

$$\leq r(s, \pi'(s)) + \gamma \sum_{s' \in \mathcal{S}} T(s, \pi'(s), s') \left[r(s', \pi'(s)) + \gamma \sum_{s'' \in \mathcal{S}} \dots \right]$$

$$\leq \dots$$

$$\leq V^{\pi'}(s).$$

Policy Iteration Algorithm always converges to π^*

If improvement stops, *i.e.* $\max_{a \in \mathcal{A}} Q^\pi(s, a) = Q^\pi(s, \pi(s)) = V^\pi(s)$, $\forall s \in \mathcal{S}$. Thus V^π satisfies the DP equation $V^\pi = \max_{a \in \mathcal{A}} Q^\pi(s, a)$, and then

$$V^\pi = V^*.$$

That means that π is optimal!

Learning and Planning

Two fundamental problems in sequential decision making under uncertainty:

- Reinforcement Learning:

- The environment is initially unknown
- The agent interacts with the environment
- The agent improves its policy

- Planning:

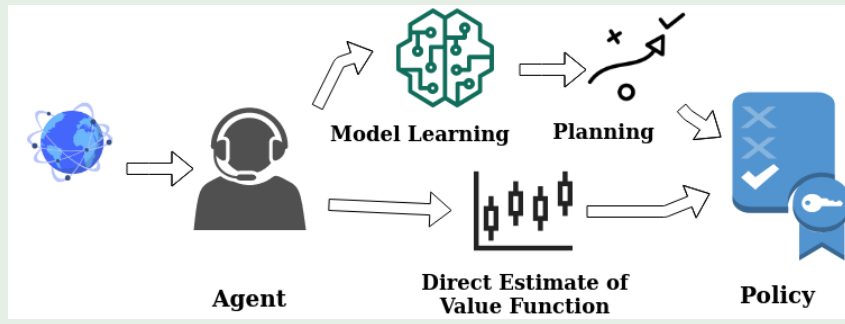
- A model of the environment is known (\mathcal{M})
- The agent performs computations with its model (without any external interaction)
- The agent improves its policy

- 1 Introduction
- 2 Markov Decision Processes
- 3 Reinforcement Learning**

Model based/Model free Reinforcement Learning

In RL, the MDP \mathcal{M} is unknown. However the agent can interact with the associated environment, and compute some estimates to improve the strategy.

- model based: estimation of T and optimization of the resulting MDP.
- model free: estimation of Q^* , or directly π^* .



Vocabulary: different types of RL

- Value Based
 - No Policy (Implicit, $\arg \max_{a \in \mathcal{A}} Q(s, a)$)
 - Value Function
- Policy Based
 - Policy
 - No Value Function
- Actor Critic
 - Policy
 - Value Function

A first learning method: averaging sample returns,
a.k.a. **Monte Carlo** method.

Monte Carlo Policy Evaluation

- Consider a trajectory using strategy π : $(S_0, S_1, \dots, S_T) \sim T^\pi$,
- Define the **return** of the state S_t :

$$G_t = \sum_{t'=t}^T \gamma^{t'-t} r(S_{t'}, \pi(S_{t'})) = R_t + \gamma R_{t+1} + \gamma^2 R_{t+1} + \dots + \gamma^{T-t} R_T,$$

with $R_t = r(S_t, \pi(S_t))$.

- Since γ^t goes to 0 exponentially, V^π is close to the expected return:

$$V^\pi(s) \approx \mathbb{E}[G_t \mid S_t = s].$$

- Monte-Carlo policy evaluation uses empirical mean return instead of expected return

Remark

The mean μ_k of a sequence $(x_i)_{i=1}^k$ can be computed incrementally:

$$\begin{aligned}\mu_{k+1} &= \frac{1}{k+1} \sum_{i=1}^{k+1} x_i = \frac{x_{k+1}}{k+1} + \frac{1}{k+1} \sum_{i=1}^k x_i \\ &= \frac{x_{k+1}}{k+1} + \frac{k}{k+1} \mu_k \\ &= \frac{1}{k+1} (x_{k+1} + k\mu_k) \\ &= \frac{1}{k+1} (x_{k+1} + (k+1)\mu_k - \mu_k) \\ &= \mu_k + \frac{1}{k+1} (x_{k+1} - \mu_k).\end{aligned}$$

First-visit Monte-Carlo prediction: estimation of V^π

input : strategy π , number of episodes n

output: V , estimation of V^π

$N(s) \leftarrow 0, \forall s \in \mathcal{S};$

for i **in** $\{1, \dots, n\}$ **do**

 generate (s_0, s_1, \dots, s_T) using π ;

 visited $\leftarrow \emptyset$;

for t **in** $\{0, \dots, T\}$ **do**

if $s_t \notin \text{visited}$ **then**

$N(s_t) \leftarrow N(s_t) + 1$;

$G(s_t) \leftarrow \sum_{i=t}^T \gamma^{i-t} r(s_i, \pi(s_i))$;

$V(s_t) \leftarrow V(s_t) + \frac{1}{N(s_t)} (G(s_t) - V(s_t))$;

 visited.add(s_t);

Remarks

- $V(s)$ is the mean of independent identically distributed (i.i.d.) samples G .
- Strong Law of Large Numbers (SLLN) $\implies V(s) \xrightarrow[n \rightarrow +\infty]{a.s.} V^\pi$.

$$\mathbb{E}[V(s)] = \mathbb{E}\left[\frac{1}{n} \sum_{i=1}^n G_i\right] = \frac{1}{n} \sum_{i=1}^n \mathbb{E}[G_i] = V^\pi(s)$$

$$\text{Var}(V(s)) = \text{Var}\left(\frac{1}{n} \sum_{i=1}^n G_i\right) = \frac{1}{n^2} \sum_{i=1}^n \text{Var}(G_i) = \frac{\text{Var}(G)}{n}$$

$$\mathbb{E}\left[\left(V(s) - V^\pi(s)\right)^2\right] = \mathbb{E}[V(s) - \mathbb{E}[V(s)]]^2 + \text{Var}(V(s)) = 0 + \frac{\text{Var}(G)}{n} \sim \frac{1}{n}.$$

- “Every visit MC” consists in removing the condition “ $s_t \notin \text{visited}$ ”.
- In this case, the samples G are not independent anymore, but it converges too.
- Q^π can be estimated in the same way.

Introduction to Temporal Differences (TD)

First, consider again the estimation update: $V(s) \leftarrow V(s) + \frac{1}{N(s)} (G(s) - V(s))$.

Fixing $N \Rightarrow$ old episodes forgotten with updates.

With $\forall s \in \mathcal{S}, N(s) = \frac{1}{\alpha} \in]0, 1[$, we get: $V(s) \leftarrow V(s) + \alpha (G(s) - V(s))$.

Another idea is to use $R_t + \gamma V(S_{t+1})$ instead of $G(S_t)$:

$$V(S_t) \leftarrow V(S_t) + \alpha (R_t + \gamma V(S_{t+1}) - V(S_t)).$$

TD vocabulary

- $R_t + \gamma V(S_{t+1})$ is called the **TD target**,
- $R_t + \gamma V(S_{t+1}) - V(S_t)$ the **TD error**.

- TD learns online after every step
→ learns from incomplete sequences.
- MC must wait until end of episode before return is known
→ learns from complete sequences.
- TD works in continuing (non-terminating) environments
- MC only works for episodic (terminating) environments

Let's consider a batch of K episodes:

$$\begin{aligned} e_1 &= s_0^{(1)}, a_0^{(1)}, r_0^{(1)}, \quad s_1^{(1)}, a_1^{(1)}, r_1^{(1)}, \quad \dots, s_{H_1}^{(1)} \\ e_2 &= s_0^{(2)}, a_0^{(2)}, r_0^{(2)}, \quad s_1^{(2)}, a_1^{(2)}, r_1^{(2)}, \quad \dots, s_{H_2}^{(2)} \\ &\vdots \\ e_K &= s_0^{(K)}, a_0^{(K)}, r_0^{(K)}, \quad s_1^{(K)}, a_1^{(K)}, r_1^{(K)}, \quad \dots, s_{H_K}^{(K)} \end{aligned}$$

Monte-Carlo:

- $K(s) = \sum_{k=1}^K \mathbb{1}_{\{s \in e_k\}} \leq K$, $t_k(s) = \min \{ t \in \{0, \dots, H_k\} \mid s_t \in e_k \text{ and } s_t = s \}$.
- $\forall s \in \mathcal{S}$ such that $K(s) > 0$, $V(s) = \frac{1}{K(s)} \sum_{\substack{k \text{ s.t.} \\ s \in e_k}} \sum_{t=t_k(s)}^{H_k-1} \gamma^{t-t_k(s)} r_t^{(k)} = \frac{1}{K(s)} \sum_{\substack{k \text{ s.t.} \\ s \in e_k}} g^{(k)}(s)$.
- $V(s)$ minimizes the least squares $\sum_{\substack{k \text{ s.t.} \\ s \in e_k}} \left(V(s) - g^{(k)}(s) \right)^2$.

Let's consider a batch of K episodes:

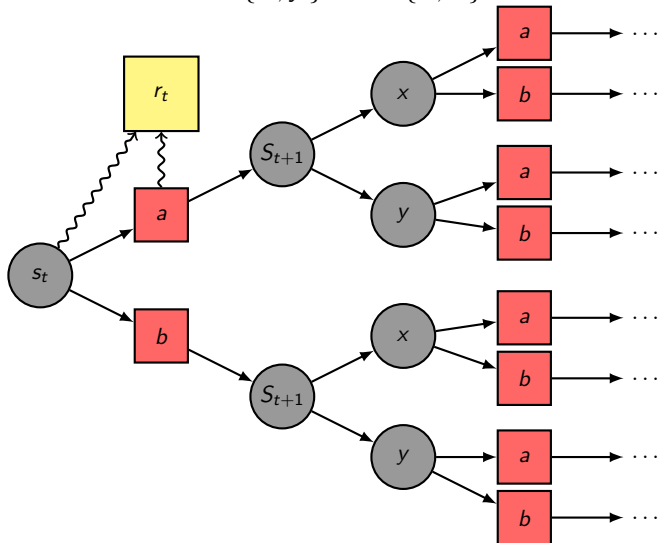
$$\begin{aligned} e_1 &= s_0^{(1)}, a_0^{(1)}, r_0^{(1)}, \quad s_1^{(1)}, a_1^{(1)}, r_1^{(1)}, \quad \dots, \quad s_{H_1}^{(1)} \\ e_2 &= s_0^{(2)}, a_0^{(2)}, r_0^{(2)}, \quad s_1^{(2)}, a_1^{(2)}, r_1^{(1)}, \quad \dots, \quad s_{H_2}^{(2)} \\ &\vdots \\ e_K &= s_0^{(K)}, a_0^{(K)}, r_0^{(K)}, \quad s_1^{(K)}, a_1^{(K)}, r_1^{(K)}, \quad \dots, \quad s_{H_K}^{(K)} \end{aligned}$$

Temporal Differences TD(0):

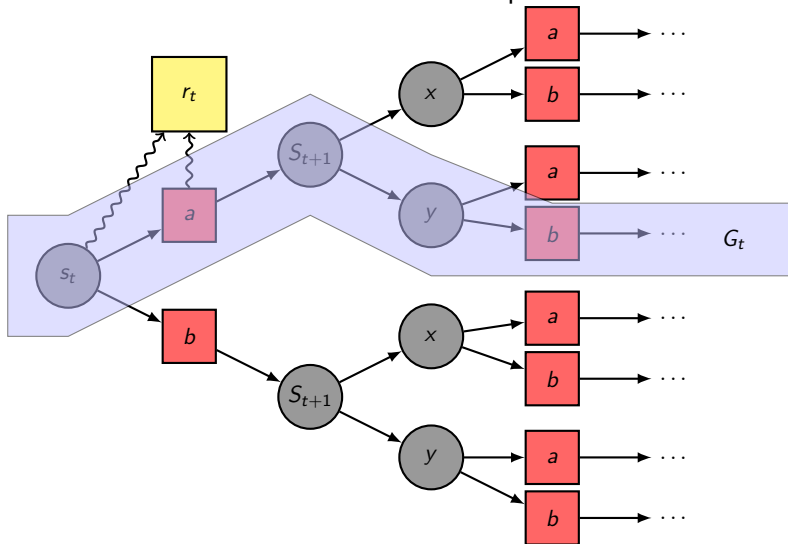
- Under some conditions, TD(0) converges to solution of max likelihood Markov model.
- Solution to the MDP $\langle \mathcal{S}, \mathcal{A}, \hat{T}, \hat{r}, \gamma \rangle$ with

- $\hat{T}(s, a, s') = \frac{1}{N(s, a)} \sum_{k=1}^K \sum_{t=0}^{H_k-1} \mathbb{1}_{\{(s, a, s') = (s_t^{(k)}, a_t^{(k)}, s_{t+1}^{(k)})\}}, \text{ and}$
- $\hat{R}(s, a) = \frac{1}{N(s, a)} \sum_{k=1}^K \sum_{t=0}^{H_k-1} r_t^{(k)} \mathbb{1}_{\{(s, a) = (s_t^{(k)}, a_t^{(k)})\}}.$

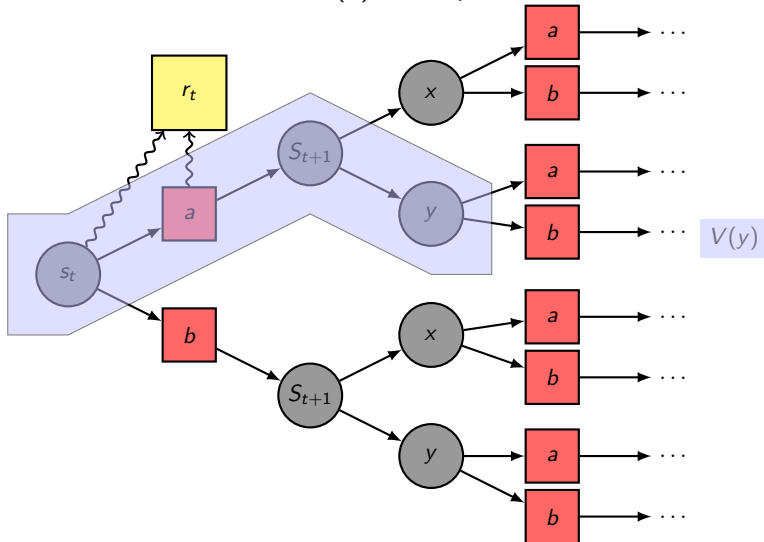
$\mathcal{S} = \{x, y\}$, $\mathcal{A} = \{a, b\}$:



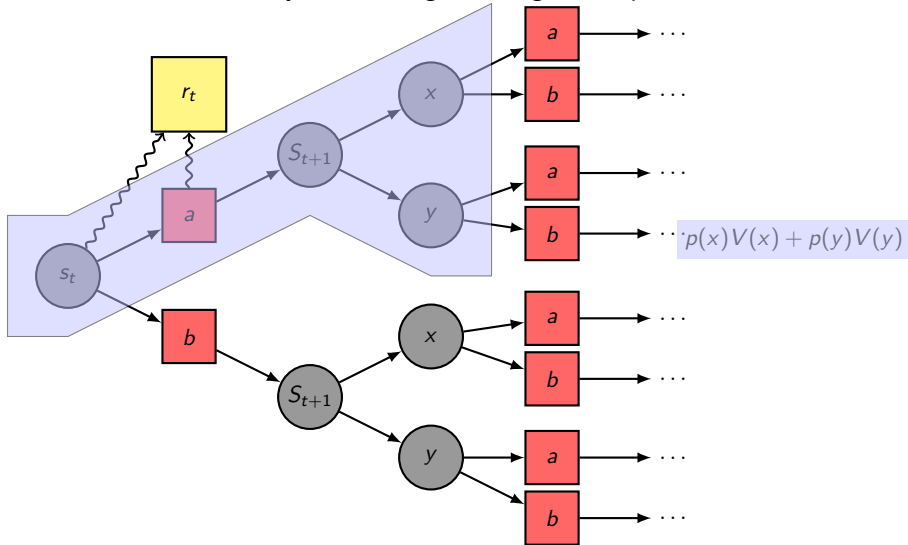
Monte Carlo Backup:



TD(0) Backup:



Dynamic Programming Backup:



Bootstrap & Sampling

- Bootstrapping: update involves an estimate
 - MC does not bootstrap
 - DP bootstraps
 - TD bootstraps
- Sampling: update samples an expectation
 - MC samples
 - DP does not sample
 - TD samples

TD(λ) with $\lambda \in [0, 1]$

- n -step return: $G_t^n = R_t + \gamma R_{t+1} + \dots + \gamma^{n-1} R_{t+n} + \gamma^n V(S_{t+n})$

- G_t^∞ is the classical return $G_t = \sum_{t'=t}^{H-1} \gamma^{t'-t} R_{t'}$

- G_t^1 is the TD target $G_t^0 = R_t + \gamma V(S_{t+1})$.

- λ -return $G_t^\lambda = (1 - \lambda) \sum_{n=1}^{\infty} \lambda^{n-1} G_t^n$

average of $(G_t^n)_{n \geq 1}$ weighted by a geometric distribution with parameter $1 - \lambda$.

- TD(λ) Backup:

$$V(S_t) \leftarrow V(S_t) + \alpha (G_t^\lambda - V(S_t)).$$

Greedy policy ?

- Greedy policy improvement over value function V **requires MDP model**:

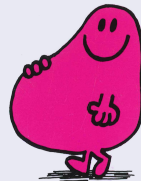
$$\pi'(s) = \arg \max_{a \in \mathcal{A}} \left\{ r(s, a) + \gamma \sum_{s' \in \mathcal{S}} T(s, a, s') V(s') \right\}.$$

- Greedy policy improvement over Q-value function Q is **model free**:

$$\pi'(s) = \arg \max_{a \in \mathcal{A}} \{ Q(s, a) \}.$$

MR. GREEDY

by Roger Hargreaves



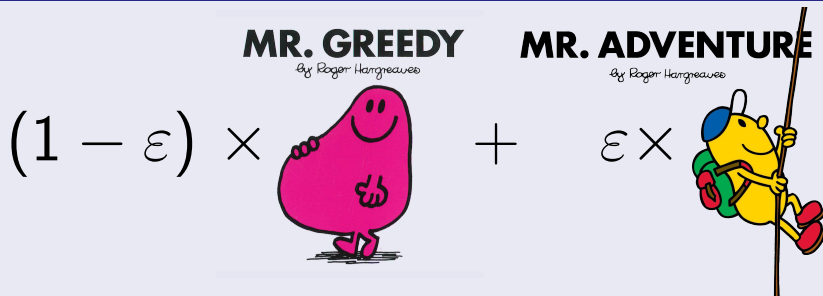
Greedy policy improvement

Let $\pi : \mathcal{S} \rightarrow \mathcal{A}$ be a strategy. Consider the greedy policy π' with respect to Q^π :

$$\pi'(s) = \arg \max_{a \in \mathcal{A}} Q^\pi(s, a).$$

Then $V^{\pi'} \geq V^\pi$ (see slide 35).

Epsilon-greedy policy ?



- For a **random policy** π , $\forall s \in \mathcal{S}$, $\pi(s)$ is a probability distribution over actions \mathcal{A} , with values denoted by $\pi(a | s) \in [0, 1]$, $\forall a \in \mathcal{A}$.
- An ϵ -**greedy** policy π is the following random policy:
$$\pi(a | s) = \begin{cases} 1 - \epsilon + \frac{\epsilon}{\#\mathcal{A}} & \text{if } a = a^*, \\ \frac{\epsilon}{\#\mathcal{A}} & \text{otherwise.} \end{cases}, \text{ with } a^* \in \arg \max_{a \in \mathcal{A}} Q(s, a) \text{ (greedy action).}$$

Epsilon-greedy policy improvement

Let $\pi : \mathcal{S} \rightarrow \mathcal{A}$ be an ε -greedy strategy. Consider the ε -greedy policy π' with respect to Q^π :

$$\pi'(a | s) = \begin{cases} 1 - \varepsilon + \frac{\varepsilon}{\#\mathcal{A}} & \text{if } a = a^*, \\ \frac{\varepsilon}{\#\mathcal{A}} & \text{otherwise.} \end{cases} \quad \text{with } a^* \in \arg \max_{a \in \mathcal{A}} Q^\pi(s, a).$$

Note that $\frac{\varepsilon}{\#\mathcal{A}} \leq \pi(a | s) \leq 1 - \varepsilon + \frac{\varepsilon}{\#\mathcal{A}}$. We have:

$$\begin{aligned} V^\pi(s) &= \sum_{a \in \mathcal{A}} \pi(a | s) Q^\pi(s, a) = (1 - \varepsilon) \sum_{a \in \mathcal{A}} \frac{\pi(a | s) - \frac{\varepsilon}{\#\mathcal{A}}}{1 - \varepsilon} Q^\pi(s, a) + \frac{\varepsilon}{\#\mathcal{A}} \sum_{a \in \mathcal{A}} Q^\pi(s, a) \\ &\leq (1 - \varepsilon) \max_{a \in \mathcal{A}} Q^\pi(s, a) + \frac{\varepsilon}{\#\mathcal{A}} \sum_{a \in \mathcal{A}} Q^\pi(s, a) \\ &= \sum_{a \in \mathcal{A}} \pi'(a | s) Q^\pi(s, a) = Q^\pi(s, \pi'(s)) \leq V^{\pi'}(s), \end{aligned}$$

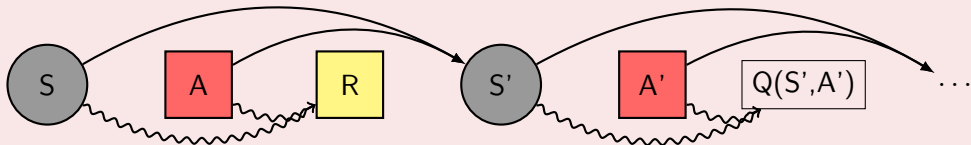
with last inequality from slide 35.

Learn and act with ϵ -greedy policies

- ϵ -greedy policies can be used with MC and TD(λ) for control
- $(1 - \epsilon)$ is the proportion of **exploitation**, ϵ of **exploration**
- another algorithm based on the estimation of the Q -value function is called SARSA.



SARSA



$$Q(s, a) \leftarrow Q(s, a) + \alpha \left(r_t + \gamma Q(s_{t+1}, a_{t+1}) - Q(s, a) \right).$$

SARSA

input : arbitrarily initialized Q-value

output: improved Q

for i **in** $\{1, \dots, n\}$ **do**

 initialize $s \sim p_{s_0}$;

$a \leftarrow \text{epsilon_greedy}(Q, s)$;

while *episode not finished* **do**

$(r, s') \leftarrow \text{next_state}(s, a)$;

$a' \leftarrow \text{epsilon_greedy}(Q, s')$;

$Q(s, a) \leftarrow Q(s, a) + \alpha(r + \gamma Q(s', a') - Q(s, a))$;

$a \leftarrow a'$;

$s \leftarrow s'$;

Remarks

- SARSA(λ) is the version with the averaged n -step returns *cf.* slide 54 for TD(λ).
- SARSA is said **on policy** because the same policy is used for control and Q -update.

What if π does not accomplish to remove the condition $\exists s : s \notin \text{visited}$?

Then we can not fix N to forget old episodes, because some $N(s) = 0 \implies \alpha = \frac{1}{N} \rightarrow +\infty$.

Idea: change π while collecting (s_0, s_1, \dots, s_T) to ensure that all "relevant" s are explored.

Set the learning rate $0 < \alpha < 1$. At each step t

Q-learning

- ε -greedy exploration strategy.

Select the action:

- random action $a \in \mathcal{A}$ with probability ε (exploration),
- $a_t \in \arg \max_{a \in \mathcal{A}} Q(s_{t+1}, a)$ with probability $1 - \varepsilon$ (exploitation).
- Receive the state the reward $r(s_t, a_t)$ and next state s_{t+1} .
- Update the Q-value estimate:

$$Q(s_t, a_t) = (1 - \alpha) \cdot Q(s_t, a_t) + \alpha \cdot \left(r(s_t, a_t) + \gamma \max_{a \in \mathcal{A}} Q(s_{t+1}, a) \right).$$

Q-learning

input : arbitrarily initialized Q-value

output: improved Q

for i **in** $\{1, \dots, n\}$ **do**

 initialize $s \sim p_{s_0}$;

$a \leftarrow \text{epsilon_greedy}(Q, s)$;

while *episode not finished* **do**

$(r, s') \leftarrow \text{next_state}(s, a)$;

$Q(s, a) \leftarrow Q(s, a) + \alpha(r + \gamma \max_{a' \in \mathcal{A}} Q(s', a') - Q(s, a))$;

$a \leftarrow \text{epsilon_greedy}(Q, s')$;

$s \leftarrow s'$;

Convergence of Q-learning

Since

- 1 in the rhs we take $\max_{a \in \mathcal{A}} Q(s_{t+1}, a)$;
- 2 ϵ -greedy exploration guarantees to explore all the required regions of $(s, a) \in \mathcal{S} \times \mathcal{A}$;
- 3 $\alpha \in]0, 1[$

\implies the estimate of Q obtained by Q-learning converges to the optimal Q^*

Tutorial Q-learning

Let's use Q-learning!

Deep Q Network

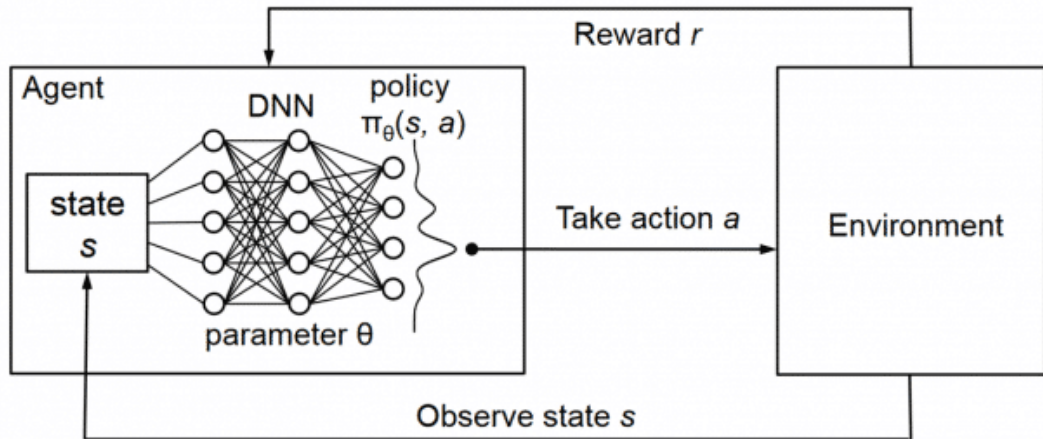
Q-learning \leftarrow unsupervised estimation of $Q(s, a)$

Starts at random \rightarrow iterates \rightarrow converges to the true $Q^*(s, a)$

What if $|\mathcal{S}| = +\infty$, as for example when $s \in \mathcal{S}$ can be described by a collection of *continuous* features? Q-learning stores a big table for each state... now it is not possible! e.g. Imagine that every state a plane can assume is represented by its position in \mathbb{R}^3 .

Idea! Instead of storing the estimates of $Q(s, a)$ in a table we can learn directly a *function approximator* that given s computes $Q(s, a) \forall a$.

This can be done in continuous state spaces using a deep neural network! [MKS⁺15]



Training of a DQN

- Least Square regression of the TD target updating the weights θ :
$$[Q_{\theta_{i+1}}(s_t, a_t) - (r_t + \gamma \max_{a'} Q_{\theta_i}(s_{t+1}, a'))]^2 \approx 0$$

where i is an iteration index of the algorithm and t a time index along a trajectory, $i \neq t$.
- Neural Networks learn fast using batches and gradient descent. The α of Q -learning is now implemented as the learning rate of the Stochastic Gradient Descent optimizer!
- In order to use batches record (s_t, a_t, r_t, s_{t+1}) in a memory called the *Replay Buffer*.
- LS regression over a batch can be done by minimizing the MSE, however the samples must be i.i.d. This is not the case along a trajectory \implies sample at random from the Replay Buffer, if it is big enough we will likely sample at each iteration i a batch containing uncorrelated transitions.

Convergence of DQN

Convergence of DQN is not guaranteed since DQN is a function approximator, but it usually happens if we use particular values for the learning rate α and the ε of an ε -greedy exploration strategy. It can be useful to define an exploration strategy with a decaying ε that starts very high $\varepsilon \approx 0.99$, so that the agent starts to explore the environment and little by little it will automatically focus on regions that are relevant for its learning of a good policy. More generally, [SB18] speaks of the *Deadly Triad* of Reinforcement Learning:

- Function Approximator
- Bootstrapping
- Off-policy strategy

When these three factors coexist in a RL scheme its convergence can be tricky.

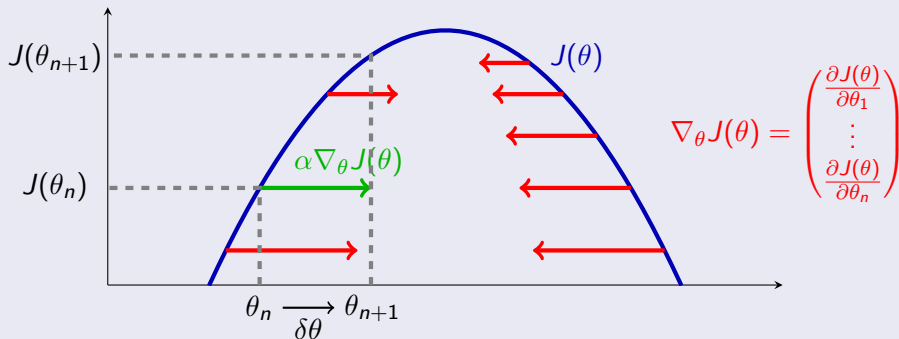
Power comes at a price.

And why not look directly for the strategy?

- Consider a random policy $\pi_\theta : \mathcal{S} \rightarrow [0, 1]^{\mathcal{A}}$. $\forall s \in \mathcal{S}, \sum_{a \in \mathcal{A}} \pi_\theta(a | s) = 1$.
- This policy is parameterized by $\theta \in \mathbb{R}^p$:
the goal is then to find θ^* such that $V^{\pi_{\theta^*}}(s)$ is high $\forall s \in \mathcal{S}$.
- We then define a score $J(\theta) = \mathbb{E}[V^{\pi_\theta}(S)]$, or $J(\theta) = \sum_{s \in \mathcal{S}} d^{\pi_\theta}(s) V^{\pi_\theta}(s)$,
with d^{π_θ} the *stationary distribution* of the Markov Chain using π_θ
(proportion of time spent in the state).
- We look for θ that maximizes $J(\theta)$: optimization problem!
- Use for example gradient descent.

In fact, it's a gradient ascent

The parameter update is $\delta\theta = \alpha \nabla_{\theta} J(\theta)$,
with $\alpha > 0$ the step-size (or learning rate) parameter, and $\nabla_{\theta} J(\theta)$ the gradient of the score.



J strictly concave \Rightarrow returns θ^* close to the global minimum *cf.* Section 5.6.2.1 of [HMS20].
 J estimated with a batch of samples \Rightarrow Stochastic Gradient Descent (SGD).

Computation of the gradient

- First idea: estimation of partial derivatives (finite differences)

$$\frac{\partial J(\theta)}{\partial \theta_k} \approx \frac{J(\theta + \varepsilon u^{(k)}) - J(\theta)}{\varepsilon},$$

where $u^{(k)} \in \mathbb{R}^p$, $u_i^{(k)} = \mathbb{1}_{\{i=k\}}$ and $\varepsilon > 0$.

- Given an estimation of $J(\theta)$, we just need an evaluation of the policy $\pi_{\theta + \varepsilon u^{(k)}}$.
- Algorithm: initialize θ and repeat until convergence
 - $\forall k \in \{1, \dots, p\}$, evaluate $\pi_{\theta + \varepsilon u^{(k)}}$ by sampling trajectories
 - use evaluation $J(\theta + \varepsilon u^{(k)})$ to compute $\nabla_{\theta} J(\theta)$
 - update parameters: $\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\theta)$

Computation of the gradient

- Another idea: if $\theta \mapsto \pi_\theta$ is differentiable.
- Let $\nabla_\theta \pi_\theta$ be its gradient.
- Let rewrite the score: $J(\theta) = \mathbb{E}[V^{\pi_\theta}(S)]$

$$= \mathbb{E}[Q^{\pi_\theta}(S, A)] = \mathbb{E}[\mathbb{E}[Q^{\pi_\theta}(S, A) \mid S]] = \mathbb{E}\left[\sum_{a \in \mathcal{A}} \pi_\theta(a \mid S) Q^{\pi_\theta}(S, a)\right].$$

- Its gradient becomes: $\nabla_\theta J(\theta) = \mathbb{E}\left[\sum_{a \in \mathcal{A}} \nabla_\theta \pi_\theta(a \mid S) Q^{\pi_\theta}(S, a)\right] =$

$$\mathbb{E}\left[\sum_{a \in \mathcal{A}} \pi_\theta(a \mid S) \frac{\nabla_\theta \pi_\theta(a \mid S)}{\pi_\theta(a \mid S)} Q^{\pi_\theta}(S, a)\right] = \mathbb{E}\left[\sum_{a \in \mathcal{A}} \pi_\theta(a \mid S) \nabla_\theta \log(\pi_\theta(a \mid S)) Q^{\pi_\theta}(S, a)\right] \\ = \mathbb{E}\left[\mathbb{E}\left[\nabla_\theta \log(\pi_\theta(A \mid S)) Q^{\pi_\theta}(S, A) \mid S\right]\right] = \mathbb{E}\left[\nabla_\theta \log(\pi_\theta(A \mid S)) Q^{\pi_\theta}(S, A)\right].$$

Monte-Carlo Policy Gradient (REINFORCE)

By using the return $G(S_t, A_t) = R_t + \gamma R_{t+1} + \dots + \gamma^{T-t} R_T$ as estimate for $Q(S_t, A_t)$,

$$\nabla_{\theta} J(\theta) \approx \nabla_{\theta} \log \left(\pi_{\theta}(a | s) \right) G(s, a).$$

input : step size α , number of episodes n

output: optimized θ

initialise θ arbitrarily;

for i **in** $\{1, \dots, n\}$ **do**

 generate (s_0, s_1, \dots, s_T) using π_{θ} ;

for t **in** $\{0, \dots, T\}$ **do**

$N(s_t) \leftarrow N(s_t) + 1$;

$G(s_t, a_t) \leftarrow \sum_{i=t}^T \gamma^{i-t} r(s_i, \pi(s_i))$;

$\theta \leftarrow \theta + \alpha \nabla_{\theta} \log \left(\pi_{\theta}(a_t | s_t) \right) G(s_t, a_t)$;

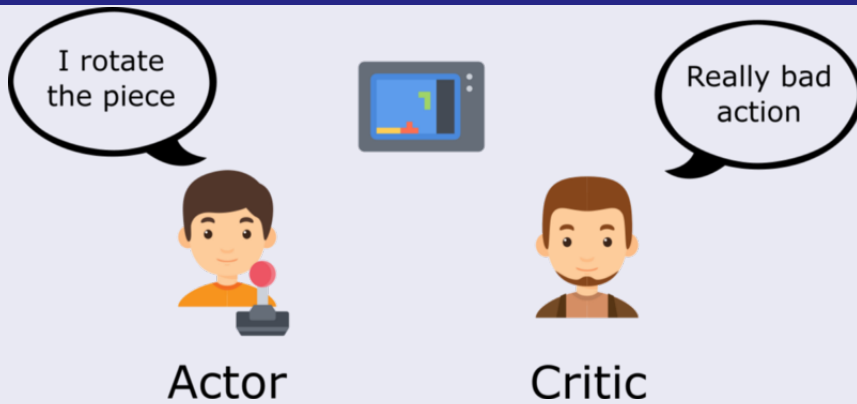
Reducing variance using a critic

- Monte-Carlo policy gradient still has high variance
- We use a **critic** to estimate the action-value function $Q_\eta \approx Q^{\pi_\theta}$, with $\eta \in \mathbb{R}^q$ containing other parameters
- Actor-critic algorithms maintain two sets of parameters
 - Critic updates action-value function parameters η
 - Actor updates policy parameters θ , in direction suggested by critic
- Actor-critic algorithms follow an approximate policy gradient

$$\begin{aligned}\nabla_\theta J(\theta) &\approx \mathbb{E}_{\pi_\theta} \left[\nabla_\theta \log \left(\pi_\theta(a | s) Q_\eta(s, a) \right) \right]. \\ \delta\theta &= \alpha \nabla_\theta \log \left(\pi_\theta(a | s) \right) Q_\eta(s, a).\end{aligned}$$

See also this introduction to Advantage Actor Critic.

Actor Critic



Institut Supérieur de l'Aéronautique et de l'Espace

10 avenue Édouard Belin – BP 54032

31055 Toulouse Cedex 4 – France

Phone: +33 5 61 33 80 80

www.isae-superaero.fr