

Reinforcement Learning

Correction TD 1 - MDP

Fabien Pesquerel
fabien.pesquerel@inria.fr

Odalric-Ambrym Maillard
odalric.maillard@inria.fr

Patrick Saux
patrick.saux@inria.fr

March 16, 2021

1 Questions from the audience

You can write below the questions that were asked during the session and see later if you can remember/re-derive the answers.

2 About the course

Exercise 1. *Given a MDP with bounded rewards, $M = (S, A, p, r)$, and a deterministic policy π ,*

- 1. Recall the Bellman operator associated to the policy π and discounted reward γ ,*
- 2. Prove that the Bellman operator is a γ -contraction for the infinity norm.*

Solution 1. Let $\pi: \mathcal{S} \rightarrow \mathcal{A}$ a (stationary) policy and $\mu_\pi(s) = \mathbb{E}_\pi[r(s, \pi(s))]$ the average reward when taking action $\pi(s)$ in state $s \in \mathcal{S}$.

1. The value of policy π at state $s \in \mathcal{S}$ is

$$\begin{aligned}
V^\pi(s) &= \mathbb{E}_\pi \left[\sum_{t \in \mathbb{N}} \gamma^t r(s_t, \pi(x_t)) \mid s_0 = s \right] \\
&= \mathbb{E}_\pi[r(s_0, \pi(s_0)) \mid s_0 = s] + \mathbb{E}_\pi \left[\sum_{t \geq 1} \gamma^t r(s_t, \pi(x_t)) \mid s_0 = s \right] \\
&= \mu_\pi(s) + \gamma \mathbb{E}_\pi \left[\sum_{t \geq 1} \gamma^{t-1} r(s_t, \pi(x_t)) \mid s_0 = s \right] \\
&= \mu_\pi(s) + \gamma \sum_{s' \in \mathcal{S}} \mathbb{P}(s_1 = s' \mid s_0 = s; \pi(s_0)) \mathbb{E}_\pi \left[\sum_{t \geq 1} \gamma^{t-1} r(s_t, \pi(x_t)) \mid s_0 = s, s_1 = s' \right] \\
&= \mu_\pi(s) + \gamma \sum_{s' \in \mathcal{S}} \mathbb{P}(s_1 = s' \mid s_0 = s; \pi(s_0)) \underbrace{\mathbb{E}_\pi \left[\sum_{t \geq 1} \gamma^{t-1} r(s_t, \pi(x_t)) \mid s_1 = s' \right]}_{=V^\pi(s')} \quad (\text{Markov property}).
\end{aligned}$$

In the tabular case ($\mathcal{S} = \{s^1, \dots, s^n\}$ is finite), one may conveniently rewrite the last equation in matrix form:

$$V^\pi = M_\pi + \gamma P_\pi V^\pi$$

where

- $V^\pi = (V^\pi(s^i))_{i=1, \dots, n}$,
- $M_\pi = (\mu_\pi(s^i))_{i=1, \dots, n}$,
- $P_\pi = \left(\mathbf{P}(s^j \mid s^i, \pi(s^i)) \right)_{i,j=1, \dots, n}$, \mathbf{P} being the MDP transition kernel.

2. The fixed point equation above calls for a Picard iteration scheme. As a reminder, the following theorem holds.

Theorem 1 (Banach-Picard fixed point theorem). *Let (\mathcal{X}, d) a complete metric space and $F: \mathcal{X} \rightarrow \mathcal{X}$ a contract mapping (K -Lipschitz for some $K < 1$).*

- (a) F admits a unique fixed point $x^* = F(x^*)$.
- (b) For any $x_0 \in \mathcal{X}$, define $(x_n)_{n \in \mathbb{N}}$ as $x_{n+1} = F(x_n)$ for $n \in \mathbb{N}$. Then

$$d(x^*, x_n) \leq K d(x^*, x_{n-1}) \leq \dots \leq K^n d(x^*, x_0).$$

In particular $\lim_{n \rightarrow +\infty} x_n = x^*$.

In most reinforcement learning settings, the value function V^π live in a Euclidean space \mathbb{R}^n or a bounded subset of a *reasonable* function space, so the assumption of a complete metric space are easy to satisfy. What is left to prove is the contraction of the linear Bellman operator:

$$\mathcal{T}_\pi: v \mapsto \mu_\pi + \gamma P_\pi v.$$

For the sake of notation, let's assume the state space \mathcal{S} is finite and V^π lives in $\mathbb{R}^{|\mathcal{S}|}$. Let $u, v \in \mathbb{R}^{|\mathcal{S}|}$ and $s \in \mathcal{S}$, we have

$$\begin{aligned}
|\mathcal{T}_\pi u(s) - \mathcal{T}_\pi v(s)| &= \gamma |P_\pi(u(s) - v(s))| \\
&= \gamma \left| \sum_{s' \in \mathcal{S}} \mathbf{P}(s' \mid s, \pi(s))(u(s') - v(s')) \right| \quad (\text{matrix product}) \\
&\leq \gamma \sum_{s' \in \mathcal{S}} \mathbf{P}(s' \mid s, \pi(s)) |u(s') - v(s')| \\
&\leq \gamma \underbrace{\sum_{s' \in \mathcal{S}} \mathbf{P}(s' \mid s, \pi(s))}_{=1} \|u - v\|_\infty \\
&= \gamma \|u - v\|_\infty.
\end{aligned}$$

Therefore, \mathcal{T}_π is γ -contractant for the topology induced by $\|\cdot\|_\infty$.

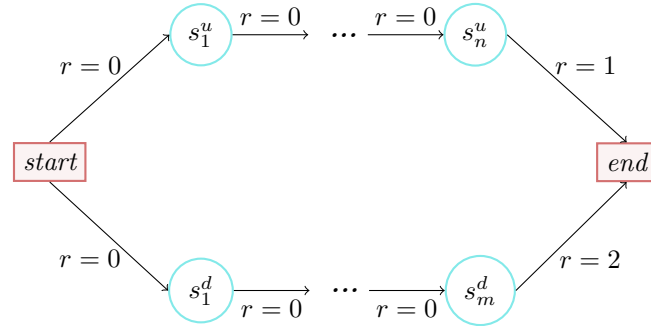
Remark 1. A similar fixed point equation with similar contraction properties holds for the optimal Bellman operator

$$\mathcal{T}: v \mapsto \max_{a \in \mathcal{A}} \mathbb{E}[r(s, a)] + \gamma P_a v$$

where $P_a = (\mathbf{P}(s' \mid s, a))_{s, s' \in \mathcal{S}}$. Contrary to \mathcal{T}_π , this operator is not linear due to the max. The proofs of fixed point and contraction are essentially the same though: simply fix an action $a \in \mathcal{A}$ instead of using $\pi(s)$, and optimise over actions at the end.

3 Exercises to build the intuition

Exercise 2 (On the influence of the discount factor). *We consider this simple MDP in which all the transitions are deterministic.*



1. Compute $Q(s, \text{up})$ and $Q(s, \text{down})$ as functions of γ (the discount factor), n (length of the upper chain, i.e., the number of states after the initial state when choosing up from that state) and m (length of the lower chain, i.e., the number of states after the initial state when choosing down from that state).
2. Compute the optimal policy π_* as a function of γ , n and m .

Solution 2. 1. There are only two policies on this deterministic MDP which can be denoted as *up* and *down*. This is because there is only one available action in each state except for the first one. Because of that, the expected values that we want to compute are just real sums. By definition of the state-action value function, one can write:

$$Q(start, up) = \mathbb{E}_{\pi_{up}, MDP} \left(\sum_{t=0}^n \gamma^t r_t | s_0 = up, a_0 = up \right) = \sum_{t=0}^n \gamma^t r_t = \gamma^n$$

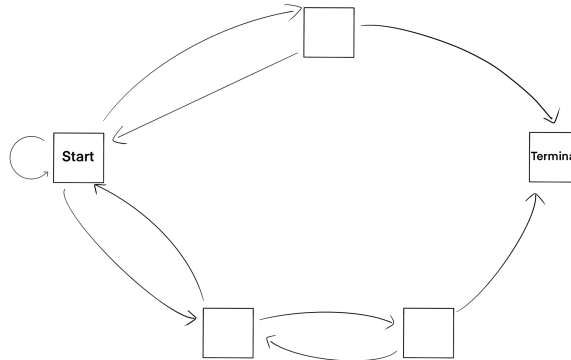
$$Q(start, down) = \mathbb{E}_{\pi_{down}, MDP} \left(\sum_{t=0}^n \gamma^t r_t | s_0 = up, a_0 = down \right) = \sum_{t=0}^n \gamma^t r_t = 2\gamma^m$$

2. Since there are only two policies, one only have to compare the two previous computed quantities:

$$\frac{Q(start, up)}{Q(start, down)} = \frac{1}{2} \gamma^{n-m}.$$

Therefore, $\pi_* = \pi_{up}$ if $\gamma^{n-m} > 2$ and $\pi_* = \pi_{down}$ if $\gamma^{n-m} < 2$. The case $\gamma^{n-m} = 2$ can be broken arbitrarily.

Exercise 3 (On the influence of the reward function). *We consider this simple maze.*



1. *Knowing that all transitions are deterministic, what are some strategies that allows to escape the maze?*
2. *Which rewards and function of rewards¹ would you choose to compute an escape policy?*
3. *Which rewards and function of rewards would you choose to compute an escape policy that uses the minimal number of moves?*

¹Discounted, average, ...

- Solution 3.** 1. We are considering deterministic policies. The space of all *valid* policies (the one that allows to reach the terminal) is made of all the policies that create a path from *start* to *terminal*. One can map *Start* to either *up* or *down*. If *up* has been chosen, one can prescribe whatever we want for the lower chain and necessarily choose the *right* action after the *up* transition. If *down* has been chosen as the first action, one can prescribe whatever we want for the upper chain and necessarily choose the *right* actions after the *up* transition in all subsequent states of the lower chain.
2. No strictly positive reward for every transitions except the ones pointing to the terminal is a necessary but not sufficient condition. For instance, if all rewards are 0s, we may not learn to escape the maze. Instead, the difference between the rewards of transitions to the terminal and other rewards should be strictly positive.
3. We may do one of the two things : modify the reward function by using a discount factor or modify the rewards of the MDP. If all rewards are 0s except the ones associated to transitions to the terminal that we set equal to 1, then the cumulative reward computed using a discount factor of γ is γ^T . T is the time spent within the maze before escaping it using a policy. The larger the time, the smaller the cumulative reward. Hence, we will compute a policy escaping with a minimal number of moves. One can also set all rewards of non-terminal transitions to -1 and set the rewards terminal transitions to 0. In this case, the cumulative reward is $-T$ where T is the same quantity as before. Once again, the smaller the time spent in the maze, the larger the cumulative reward.

4 Apply the theory

Students can hope to come back to this section later in the course and test what they will learn (algorithms, methods, heuristics). In particular, Deep RL ideas might be tested using the following exercise. Students are assumed to be familiar with *python 3.X* and usual scientific libraries. In particular, it is recommended to install *numpy*, *scipy*², *gym* and *pyTorch* (that will be used later in the course).

Exercise 4 (Gym & Frozen Lake). *Frozen Lake* is a *gym* environment that we will use to implement some algorithms that were learned during the first class. In the following, we will use a discount factor of $\gamma = 0.9$.

1. To check that everything is installed on your computer/environment run the small script [check_env.py](#). It will try to import the necessary libraries and print current versions of them. You can use those versions to help you with any troubleshooting.
2. To better grasp the Frozen Lake environment, read, understand and run [discover.py](#) (until the stop mark). Describe the MDP associated to this environment and formalize mathematically the goal of an agent.
3. Using the aforementioned environment, get a Monte-Carlo estimation of the value function of a simple deterministic policy at s_0 , the initial state. For instance, this simple policy could be to always going to the right. This is the implemented example in [discover.py](#).
4. In the case of finite MDPs, the value function V^π of a policy π is a vector in $\mathbb{R}^{|S|}$ (Tabular case). Write V^π as the solution of a the product of a matrix with a vector (both known once π is). Write a function that takes a deterministic policy π as an input and outputs its value function V^π . Compare the result of this question with the Monte-Carlo estimation of the previous question and check for consistency.
5. Using the contraction property of the Bellman operator, implement a function that iteratively apply the Bellman operator to compute the value function V^π of a policy π . What could be a good stopping criterion for your algorithm?
6. Using the contraction property of the optimal Bellman operator, implement a function that iteratively apply the optimal Bellman operator to compute the optimal value function V^* .
7. Compute an optimal policy for this environment using Value Iteration (see appendix).
8. Compute an optimal policy for this environment using Policy Iteration (see appendix).
9. Compute an optimal policy for this environment using Q-learning (see appendix). Note the change of setting: you need to simulate several episodes of Frozen Lake and learn from observed transitions only (for instance you may not assume the transition probabilities are known). Feel free to play with the number of episodes, learning rate and exploration parameter ϵ . You may set the `is_slippery` argument to `False` when declaring the environment to have deterministic transitions and help the algorithm learn faster.
10. Compare quantitatively the three methods.

²Those that are looking for performance might want to take a look at sparse representation of matrices (`scipy.sparse.csr_matrix`) and computation's techniques with arrays (BLAS and LAPACK are used as routine in `scipy`).

A Algorithms

Algorithm: Value Iteration

Input: MDP transition kernel $\mathbf{P}(s' | s, a)$ and average reward $\mu(s, a)$, discount factor γ , number of iterations N .

Initialisation: $v_0 \in \mathbb{R}^{|S|}$, $n = 0$.

while $n < N$ **do**

for $s \in \mathcal{S}$ **do**

$v(s) \leftarrow \max_{a \in \mathcal{A}} \mu(s, a) + \gamma \sum_{s' \in \mathcal{S}} P(s' | s, a) v(s')$;

end

$n \leftarrow n + 1$;

end

Return: $\pi \in \arg \max_{a \in \mathcal{A}} \mu(\cdot, a) + \gamma \sum_{s' \in \mathcal{S}} P(s' | \cdot, a) v(s')$.

Algorithm: Policy Iteration

Input: MDP transition kernel $\mathbf{P}(s' | s, a)$ and average reward $\mu(s, a)$, discount factor γ .

Initialisation: $\pi \in \mathcal{A}^{\mathcal{S}}$.

while *policy not stable* **do**

for $s \in \mathcal{S}$ **do**

$v_{\pi}(s) \leftarrow (I - \gamma \mathbf{P}_{\pi})^{-1} \mu_{\pi}(s)$; // policy evaluation

$\pi(s) \leftarrow \arg \max_{a \in \mathcal{A}} \mu(s, a) + \gamma \sum_{s' \in \mathcal{S}} P(s' | s, a) v(s')$; // policy improvement

end

end

Return: π .

Algorithm: Q-Learning (episodic, ε greedy)

Input: Discount factor γ , learning rates $(\alpha_t)_{t \in \mathbb{N}}$, number of episodes T , $\varepsilon > 0$.

Initialisation: $Q \in \mathbb{R}^{\mathcal{S} \times \mathcal{A}}$, $t = 0$.

while $t < T$ **do**

$h \leftarrow 0$;

 Start episode t in state s_h and action a_h ;

while *episode not terminated* **do**

 Take action a_h in state s_h , observe reward r_h and move to state s_{h+1} ;

$Q(s_h, a_h) \leftarrow Q(s_h, a_h) + \alpha_t (r_h + \gamma \max_{b \in \mathcal{A}} Q(s_{h+1}, b) - Q(s_h, a_h))$;

$a_{h+1} \in \arg \max_{b \in \mathcal{A}} Q(s_{h+1}, b)$ with probability $1 - \varepsilon$ else $a_{h+1} \sim \text{Unif}(\mathcal{A})$;

 Check if episode t is terminated;

end

$t \leftarrow t + 1$;

end

Return: $\pi(s) \in \arg \max_{a \in \mathcal{A}} Q(s, a)$
