

ROB311-TD2

Iad ABDUL-RAOUF and Madeleine BECKER

September 2020, 28th

Preliminary observations

It is said in the questions of this TP2 that $\gamma \in [0, 1]$ (for instance see question 4). However the environment does not contain any terminal states and does contain only positive rewards with some rewards being strictly positive. Thus $\gamma = 1$ would lead V^* to take infinite values in the iterative value algorithm. Therefore, we assumed it was a misprint and that γ actually lies in the interval $[0, 1[$ (french notation to exclude 1 from the interval).

Moreover, in the formulas of V^* and π^* , it is our understanding that \max_a and argmax_a implicitly implies that only possible actions a are taken into account. In order to make everything explicit, let us introduce the following notation :

$$A(S) = \{a \mid a \text{ is a possible action from state } S\}$$

This notation will be used in the demonstration of the two lemmas before questions 3 and 4.

Question 1

All states only have the action a_0 associated with, except for the state S_0 which can be associated with the actions a_1 and a_2 (see figure 4).

Therefore, there are only two possible policies :

$$\begin{aligned} \pi_1 : (S_0, S_1, S_2, S_3) &\longmapsto (a_1, a_0, a_0, a_0) \\ \text{and } \pi_2 : (S_0, S_1, S_2, S_3) &\longmapsto (a_2, a_0, a_0, a_0) \end{aligned}$$

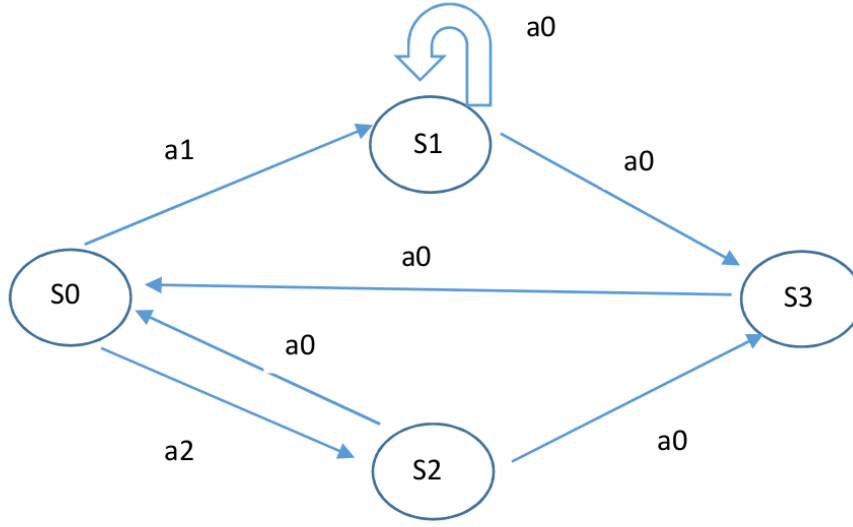


Figure 1: Graph of the links between states and associated actions

Question 2

The TD statement gives the following formula :

$$V^*(S) = R(S) + \max_a \left(\gamma \sum_{S'} T(S, a, S') V^*(S') \right) \quad (1)$$

From there we can deduce for each state :

$$\begin{aligned} V^*(S_0) &= R(S_0) + \max_a (\gamma T(S_0, a_1, S_1) V^*(S_1), \gamma T(S_0, a_2, S_2) V^*(S_2)) \\ &= 0 + \max_a (\gamma * 1 * V^*(S_1), \gamma * 1 * V^*(S_2)) \\ &= \max_a (\gamma V^*(S_1), \gamma V^*(S_2)) \end{aligned}$$

$$\begin{aligned} V^*(S_1) &= R(S_1) + \max_a (\gamma T(S_1, a_0, S_1) V^*(S_1) + \gamma T(S_1, a_0, S_3) V^*(S_3)) \\ &= \gamma(1 - x) V^*(S_1) + \gamma x V^*(S_3) \end{aligned}$$

$$\begin{aligned} V^*(S_2) &= R(S_2) + \max_a (\gamma T(S_2, a_0, S_0) V^*(S_0) + \gamma T(S_2, a_0, S_3) V^*(S_3)) \\ &= 1 + \gamma(1 - y) V^*(S_0) + \gamma y V^*(S_3) \end{aligned}$$

$$\begin{aligned} V^*(S_3) &= R(S_3) + \max_a (\gamma T(S_3, a_0, S_0) V^*(S_0)) \\ &= 10 + \gamma V^*(S_0) \end{aligned}$$

Lemmas

To answer the two following questions, we are first going to demonstrate two lemmas :

Lemma. *If all rewards are positive then $\forall S, V^*(S) \geq 0$*

and

Lemma. *For a given $\gamma \in [0, 1[$, V^* is upper-bounded.*

Demonstration of the first lemma

We have the formula (1) :

$$V^*(S) = R(S) + \max_{a \in A(S)} \left(\gamma \sum_{S'} T(S, a, S') V^*(S') \right) \quad (2)$$

where $A(S)$ is the set of possible actions from S (see Preliminary observations).

We define S_{min} , the state which minimizes $V^*(S)$. Then, for all states S :

$$\begin{aligned} V^*(S) &\geq R(S) + \max_{a \in A(S)} \left(\gamma \sum_{S'} T(S, a, S') V^*(S_{min}) \right) \\ &\geq R(S) + \max_{a \in A(S)} \left(\gamma V^*(S_{min}) \sum_{S'} T(S, a, S') \right) \\ &\geq R(S) + \max_{a \in A(S)} (\gamma V^*(S_{min}) \times 1) \\ &\geq R(S) + \gamma V^*(S_{min}) \end{aligned}$$

Thus, applied in particular to $S = S_{min}$:

$$\begin{aligned} V^*(S_{min}) &\geq R(S_{min}) + \gamma V^*(S_{min}) \\ \iff (1 - \gamma) V^*(S_{min}) &\geq R(S_{min}) \geq \min_S (R(S)) \\ \iff V^*(S_{min}) &\geq \frac{\min_S (R(S))}{1 - \gamma} \end{aligned}$$

Here are all rewards positive, therefore for all states S :

$$\boxed{V^*(S) \geq V^*(S_{min}) \geq 0}.$$

Demonstration of the second lemma

In the same way, we define S_{max} , the state which maximizes $V^*(S)$. Then, for all states S :

$$\begin{aligned} V^*(S) &\leq R(S) + \max_{a \in A(S)} \left(\gamma \sum_{S'} T(S, a, S') V^*(S_{max}) \right) \\ &\leq R(S) + \gamma V^*(S_{max}) \end{aligned}$$

Thus, applied in particular to $S = S_{max}$,

$$\begin{aligned} V^*(S_{max}) &\leq R(S_{max}) + \gamma V^*(S_{max}) \\ \iff (1 - \gamma)V^*(S_{max}) &\leq R(S_{max}) \leq \max_S(R(S)) \\ \iff V^*(S_{max}) &\leq \frac{\max_S(R(S))}{1 - \gamma} \end{aligned}$$

Yet we have $\forall S, V^*(S) \leq V^*(S_{max})$.

Therefore, for a given $\gamma \in [0, 1[$, the value function is upper-bounded.

Question 3

We are looking for a value of x for which $\forall \gamma \in [0, 1[, y \in [0, 1], \pi^*(S_0) = a_2$, where :

$$\pi^*(S) = \operatorname{argmax}_a \left(\sum_{S'} T(S, a, S') V^*(S') \right). \quad (3)$$

For $S = S_0$:

$$\pi^*(S) = \operatorname{argmax}_a (V^*(S_1), V^*(S_2))$$

Yet we have on first hand :

$$V^*(S_1) = \gamma(1 - x)V^*(S_1) + \gamma x V^*(S_3)$$

With $x = 0$:

$$\begin{aligned} V^*(S_1) &= \gamma V^*(S_1) \\ \iff (1 - \gamma)V^*(S_1) &= 0 \end{aligned}$$

With $\gamma < 1$, we can deduce that $V^*(S_1) = 0$.

On the other hand,

$$V^*(S_2) = 1 + \gamma(1 - y)V^*(S_0) + \gamma y V^*(S_3)$$

Thanks to the first lemma, we can state that $\gamma(1 - y)V^*(S_0) + \gamma y V^*(S_3) \geq 0$, therefore $V^*(S_2) \geq 1 > 0 = V^*(S_1)$. We can deduce from that :

$$\boxed{x = 0 \implies \pi^*(S_0) = a_2}$$

Question 4

We want to prove that the TD statement is wrong, therefore that $\forall \gamma \in [0, 1], \forall y \in [0, 1], \exists x > 0 \mid \pi^*(S_0) = a_2$.

As before,

$$\pi^*(S) = \operatorname{argmax}_a (V^*(S_1), V^*(S_2))$$

We already have proven at the previous question that

$$\forall x \in [0, 1], V^*(S_2) = 1 + \gamma(1 - y)V^*(S_0) + \gamma y V^*(S_3) \geq 1.$$

We want to prove that there exists an x value such as $V^*(S_1) < 1$, because then it would be lower than $V^*(S_2)$, and then $\pi^*(S_0) = a_2$.

$$V^*(S_1) = \gamma(1 - x)V^*(S_1) + \gamma x V^*(S_3) \iff V^*(S_1) = \frac{\gamma x}{1 - \gamma(1 - x)} V^*(S_3)$$

Yet, since the highest reward is 10, the second lemma states that :

$$V^*(S_1) = \frac{\gamma x}{1 - \gamma(1 - x)} V^*(S_3) \leq \frac{\gamma x}{1 - \gamma(1 - x)} \frac{10}{1 - \gamma} \xrightarrow{x \rightarrow 0} 0$$

Therefore, for a given $0 \leq \gamma < 1$, there exists an x value close enough to zero such as $V^*(S_1) < 1 = V^*(S_2)$, and for this value of x , $\pi^*(S_0) = a_2$.

Question 5

Python code

Please, find our code in the following github repository : <https://github.com/MalphaBecker/ROB311-TD2>.

It contains our main file applying the value iteration algorithm to the graph of the TP by calling our functions defined in the file `aux_function.py`.

You just have to run `main_iad_madeleine.py` in your IPython console. Parameters of the algorithm and parameters of the graph are at the beginning of the main file and can be changed at will.

Results

To measure convergence, we used the root mean square value (RMS method). Having $x = y = 0.25, \gamma = 0.9$, and the stop condition ϵ on the RMS value set to $\epsilon = 0.01$, leading to a convergence in 50 iterations, the optimal policy found by our algorithm is :

$$\pi^* : (S_0, S_1, S_2, S_3) \longmapsto (a_1, a_0, a_0, a_0)$$

You can also find the final values of V^* , the evolution of the RMS and the evolution of the policy on the figures below which clearly show convergence. For the policy, only $\pi^*(0)$ is shown as the others are constant equal to a_0 .

```

In [1]: runfile('/home/iad/Documents/ENSTA/3A/rob_311/TP_rob311/TP2/ROB311-TD2/main_i
iad/Documents/ENSTA/3A/rob_311/TP_rob311/TP2/ROB311-TD2')
the estimated optimal policy is : [1. 0. 0. 0.]
number of iteration : 50
final V = [14.1025855  15.67876763 15.61484469 22.68402157]
final RMS : 0.009228192433676341
parameters : gamma = 0.9 , eps = 0.01

```

Figure 2: Final results of the algorithm

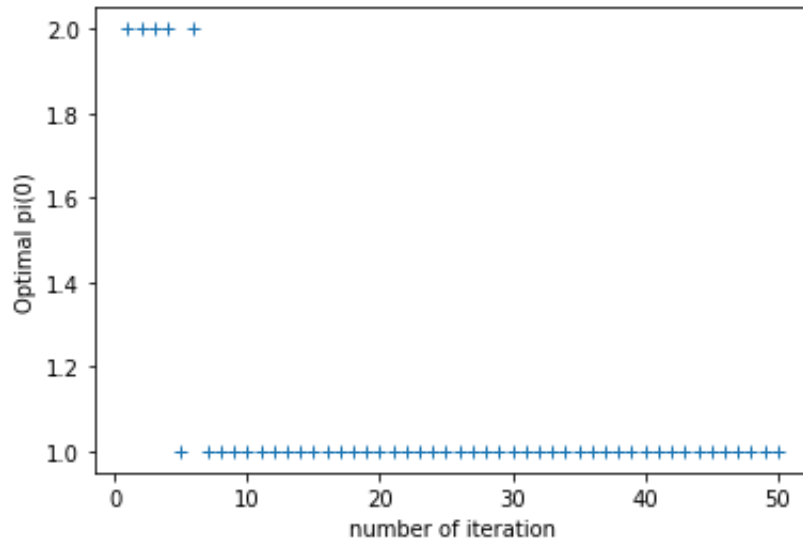


Figure 3: Optimal policy for S_0 at each iteration until convergence

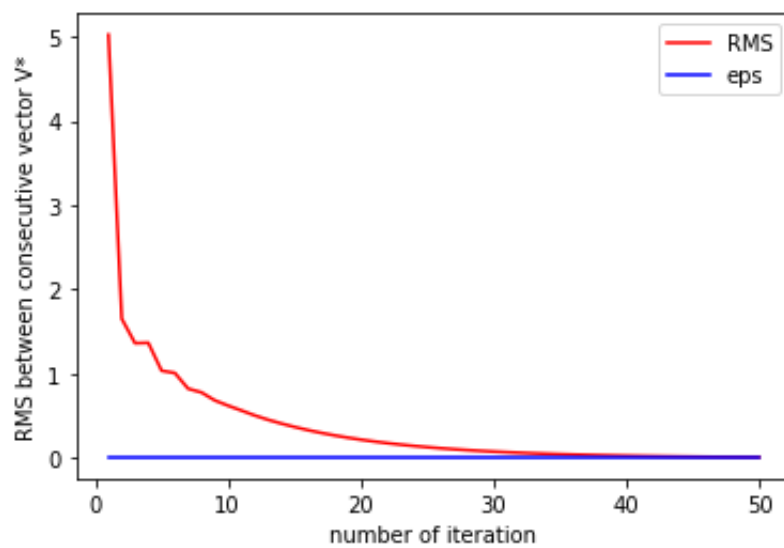


Figure 4: Root mean square error between two consecutive values of V^*