

---

Assignment 1

# Reinforcement Learning in a Discrete Domain

---

## DOMAIN

We consider the following components:

- State space:  $X = \{(x, y) \in \mathbb{N}^2 \mid x < n, y < m\}$ ,
- Action space:  $U = \{(1, 0), (-1, 0), (0, 1), (0, -1)\}$ ,
- Reward signal:  $r((x, y), (i, j)) = R(g, F((x, y), (i, j)))$ ,
- Discount factor:  $\gamma = 0.99$ ,
- Time horizon:  $T \rightarrow +\infty$ ,

where  $g \in \mathbb{R}^{n \times m}$  for any  $n, m \in \mathbb{N}$ ,  $F((x, y), (i, j)) = (\min(\max(x+i, 0), n-1), \min(\max(y+j, 0), m-1))$  and  $R(g, (x, y))$  returns the value of the cell located at  $(x, y)$  in  $g$ . Note that the reward signal is bounded by  $B = \max_{(x,y)} R(g, (x, y))$ . From these components, we derive two different domains:

- **Deterministic domain:** The dynamics  $f$  are defined as follows:

$$f((x, y), (i, j)) = F((x, y), (i, j)), \forall ((x, y), (i, j)) \in X \times U.$$

- **Stochastic domain:** The dynamics  $f$  are defined as follows:

$$f((x, y), (i, j), w) = \begin{cases} F((x, y), (i, j)) & \text{if } w \leq \frac{1}{2}, \\ (0, 0) & \text{otherwise,} \end{cases} \forall ((x, y), (i, j)) \in X \times U,$$

where  $0 \leq w \leq 1$  is drawn according to a uniform distribution. In this domain, the reward signal is redefined by  $r((x, y), (i, j), w) = R(g, f((x, y), (i, j), w))$ .

Figure 1 describes an instance of the domain that you must use in all sections of this assignment.

-3	1	-5	0	19
6	3	8	9	10
5	-8	4	1	-8
6	-9	4	19	-5
-20	-17	-4	-3	9

**Figure 1:** Domain instance. The number in each cell represents the reward obtained while reaching it after a move (e.g., if the agent moves up from the red cell, it receives a reward equal to 5).  $(0, 0)$  corresponds to the upper-left corner of the grid. The two coordinates refer to the line and to the column number, respectively. The initial state of this instance is the red cell located at  $(3, 0)$ .

## INSTRUCTIONS

You are expected to deliver (i) your source code (in Python 3.7, a file per section and named as *sectionK.py* where  $K$  is the section number - each missing one will cost you a one point penalty) and (ii) a report which is structured according to the next sections. We insist on the fact that the report is mandatory and that the source code needs to use only standard programming libraries plus, possibly, *NumPy* and *matplotlib*. We should also be able to execute your code, understand easily the results displayed. Both domains must be considered while addressing each section of this assignment. Non-deterministic numerical results should be averaged through 10 independent runs and displayed with standard deviation.

### 1 IMPLEMENTATION OF THE DOMAIN (2 POINTS)

Implement the different components of the domain instance described in Figure 1. Implement a rule-based policy of your choice (e.g., always go left, select actions always at random...) and describe it in your report. Simulate the policy in the domain through a single trajectory of 10 steps, starting by the initial state  $x_0 = (3, 0)$ . Show the trajectories in your report as a sequence of four tuples  $(x_0, u_0, r_0, x_1), \dots, (x_{10}, u_{10}, r_{10}, x_{11})$  (e.g., by a screenshot of your terminal).

### 2 EXPECTED RETURN OF A POLICY (2 POINTS)

Implement a routine to estimate  $J_\mu$  in this domain, where  $\mu : X \rightarrow U$  is a stationary policy. Test your implementation with your rule-based policy of Section 1. Choose a  $N$  which is large enough to approximate well the infinite time horizon of the policy and motivate your choice using the discount factor  $\gamma$ . Display  $J_\mu^N(x)$  for each state  $x$  in a table in your report.

### 3 OPTIMAL POLICY (3 POINTS)

Implement routines which compute  $p(x'|x, u)$  and  $r(x, u)$  that define the equivalent MDP of the domain. Afterwards, implement a routine which computes the sequence of  $Q_N$ -functions from these components using the dynamic programming principle. Determine the smallest  $N$  such that a greater value does not change the policy inferred from the last  $Q$ -function of the sequence. Derive directly  $\mu^*$  from  $Q$  and show it in a table in your report. Display  $J_{\mu^*}^N$  for each state  $x$  in a table in your report.

### 4 SYSTEM IDENTIFICATION (5 POINTS)

Implement a routine which estimates  $r(x, u)$  and  $p(x'|x, u)$  from a given trajectory  $h_t = (x_0, u_0, r_0, x_1, u_1, r_1, \dots, u_{t-1}, r_{t-1}, x_t)$ . Display the convergence speed of  $\hat{p}$  and  $\hat{r}$  towards  $p$  and  $r$  in a curve plot using the infinite norm through a growing trajectory generated by a random uniform policy. Compute  $\hat{Q}$  by using  $\hat{r}(x, u)$  and  $\hat{p}(x'|x, u)$  that estimate the MDP structure. Display, in your report, the infinite norm between  $\hat{Q}$  and  $Q$  for each trajectory length in a table in your report. Derive  $\hat{\mu}^*$  from  $\hat{Q}$  and display it in a table in your report. Display  $J_{\hat{\mu}^*}^N$ , along with  $J_{\mu^*}^N$ , for each state  $x$  in a table in your report. Explain the influence of the length of the trajectory on the quality of the approximation  $\hat{Q}$ .

### 5 Q-LEARNING IN A BATCH SETTING (8 POINTS)

#### 5.1 OFFLINE Q-LEARNING (2 POINTS)

Sample a trajectory from a random uniform policy with a finite, large enough time horizon  $T$ . Motivate your choice of  $T$ . Implement a routine which iteratively updates  $\hat{Q}(x_k, u_k)$  with  $Q$ -learning from this trajectory. Run your routine with a constant learning rate  $\alpha = 0.05$ . Derive directly  $\hat{\mu}^*$  from  $\hat{Q}$  and display it in a table in your report. Display  $J_{\hat{\mu}^*}^N$ , along with  $J_{\mu^*}^N$ , for each initial state  $x$ , in a table in your report.

#### 5.2 ONLINE Q-LEARNING (4 POINTS)

Implement an intelligent agent which learns  $Q$  with Q-learning using an  $\epsilon$ -greedy policy.

The first experimental protocol is the following. The agent trains over 100 episodes having 1000 transitions in the domain instance described in Figure 1. An episode always starts from the initial state (see Figure 1). The learning rate  $\alpha$  is equal to 0.05 and the exploration rate  $\epsilon$  is equal to 0.5. The values of  $\alpha$  and  $\epsilon$  are both constant over time. The function  $\hat{Q}$  is updated after every transition. The transitions are used only once for updating  $\hat{Q}$ .

The second experimental protocol differs from the first one only by the learning rate which is such that  $\alpha_0 = 0.05$  and  $\forall t > 0, \alpha_t = 0.8\alpha_{t-1}$ .

The third experimental protocol is identical from the first one except that the one-step system transitions are stored in a replay buffer and that at each time-step, the function  $\hat{Q}$  is updated ten times by drawing ten transitions at random from the replay buffer.

Run the three experimental protocols and display for each of them, after each episode, the value of  $\|J_{\mu_{\hat{Q}}}^N - J_{\mu^*}^N\|_{\infty}$  – where  $\mu_{\hat{Q}}$  is the policy derived from  $\hat{Q}$  – in a plot. Compare the results obtained. You might also plot  $\|J_{\mu_{\hat{Q}}}^N - J_{\mu^*}^N\|_2$  after each episode if needed.

### 5.3 DISCOUNT FACTOR (2 POINTS)

Rerun the first abovementioned experimental protocols with  $\gamma = 0.4$ . What do you observe in terms of convergence rate of the algorithm expressed as the speed at which  $\hat{Q}$  converges to the  $Q$ -function of the new domain ?

### 5.4 Q-LEARNING WITH ANOTHER EXPLORATION POLICY (BONUS, 2 POINTS)

Implement another exploration/exploitation policy presented in the scientific literature. Rerun the first experience from 6.2 by replacing the  $\epsilon$ -greedy policy with this new policy. Tune the parameters of your policy so as to maximize the expected return over a single infinite time horizon. Discuss the results obtained at the light of the exploration/exploitation tradeoff.