

---

Assignment 1

# Reinforcement Learning in a Discrete Domain

---

## 1 DOMAIN

We describe the domain below :

- State space :  $X = \{(x, y) \in \mathbb{N}^2 \mid x < n, y < m\}$ ,
- Action space :  $U = \{(1, 0), (-1, 0), (0, 1), (0, -1)\}$ ,
- Dynamics :  $(x_{t+1}, y_{t+1}) = (\min(\max(x_t + i, 0), n - 1), \min(\max(y_t + j, 0), m - 1))$   
where  $(i, j) \in U$ ,
- Reward signal :  $r((x_t, y_t), (i_t, j_t)) = R(g, x_{t+1}, y_{t+1})$ ,
- Discount factor :  $\gamma = 0.99$ ,
- Time horizon :  $T \rightarrow +\infty$ ,

where  $g \in \mathbb{R}^{n \times m}$  for any  $n, m \in \mathbb{N}$  and  $R(g, i, j)$  returns the value of the cell located at  $(i, j)$  in  $g$ . Note that the reward signal is bounded by  $B = \max_{(i,j)} R(g, i, j)$ . The domain is deterministic. We provide also a stochastic setting which differs from the deterministic setting by its dynamics, redefined below :

$$(x_{t+1}, y_{t+1}) = \begin{cases} (\min(\max(x_t + i, 0), n - 1), \min(\max(y_t + j, 0), m - 1)) & \text{if } w_t \leq 1 - \beta \\ (0, 0) & \text{otherwise,} \end{cases}$$

where  $0 \leq w_t \leq 1$  is drawn, following a uniform distribution, at each time step  $t = 0, 1, 2 \dots T$  and  $0 \leq \beta \leq 1$ .

Make sure you (i) consider both settings for each exercise, (ii) rigorously test your implementation at each stage of the assignment and (iii) store any intermediate result to avoid redundant computations. Figure 1 illustrates an instance of the domain that you will have to use for carrying out your experiments.

-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. For example, if at time  $t$  the agent moves up from the red cell, it receives a reward equal to 5.

## 2 IMPLEMENTATION OF THE DOMAIN

Implement the different components of the domain. Test your results by simulating a simple policy (e.g., always go up).

## 3 EXPECTED RETURN OF A POLICY

Implement a routine to estimate  $J^\mu$  in this domain, where  $\mu : X \rightarrow U$  is a stationary policy. The implementation needs to handle both deterministic and stochastic settings and should exploit the dynamic programming principle where a sequence of functions  $J_N$  for  $N = 0, 1, 2, \dots$  is computed using the Bellman Equation. Test your implementation with the following policy :  $\mu(x) = (0, 1), \forall x \in X$ . Display  $J_\mu^N(x)$  for each state  $x$ .

## 4 OPTIMAL POLICY

Implement a routine which computes  $p(x'|x, u)$  and  $r(x, u)$  that define the equivalent MDP of the domain, together with  $\gamma$  and  $T$ . Afterwards, implement a routine which computes  $Q(x, u)$  from these components using the dynamic programming principle. Derive directly  $\mu^*$  from  $Q$ . Display  $J_{\mu^*}^N$  for each initial state  $x$ .

## 5 SYSTEM IDENTIFICATION

Implement a routine which estimates  $r(x, u)$  and  $p(x'|x, u)$  from a given trajectory  $h_t = x_0 u_0 r_0 x_1 \dots x_{t-1} u_{t-1} r_{t-1} x_t$ . Display the convergence speed toward  $p$  and  $r$ . Compute  $\hat{Q}$  with  $\hat{r}(x, u)$  and  $\hat{p}(x'|x, u)$  the MDP estimation of the domain. Derive directly  $\hat{\mu}^*$  from  $\hat{Q}$ . Display  $J_{\hat{\mu}^*}^N$ , along with  $J_{\mu^*}^N$ , for each initial state  $x$ . Test your implementation on several trajectories. Detail your experimental protocol.

## 6 Q-LEARNING IN A BATCH SETTING

Implement a routine which computes  $\hat{Q}(x_k, u_k)$  with *Q-learning*. This routine needs to handle both single and multiple trajectories, along with the possibility to do experience replay. Display the convergence speed of  $\hat{Q}$  toward  $Q$ . Derive directly  $\hat{\mu}^*$  from  $\hat{Q}$ . Display  $J_{\hat{\mu}^*}^N$ , along with  $J_{\mu^*}^N$ , for each initial state  $x$ .

Implement an intelligent agent which learns  $Q$  with Q-learning and uses an  $\epsilon$ -greedy policy. The protocol is the following. The agent trains over 100 episodes having 1000 transitions in the domain instance described in Figure 1. An episode always start from state (3,3). The learning ratio  $\alpha_k$  is equal to 0.05,  $\forall k$ . Test your implementation with different batch sizes for experience replay.

Display  $\hat{Q}$  and  $J_{\hat{\mu}^*}^N$ , still along with  $J_{\mu^*}^N$ , with  $\hat{\mu}^*$  derived from  $\hat{\mu}^*$  under the protocol described above with  $\epsilon$  equals to 0.05, 0.2, and 0.5. Consider also a set of  $\epsilon$ 's that starts from one of these values and decreases linearly.

Implement another exploration/exploitation policy presented in the scientific literature and test it. Discuss the observed results.

## 7 DISCOUNT FACTOR

Redo the whole assignment with  $\gamma = 0.4$ . Compare the differences between the results and discuss them.