

# TD1

## Artificial Intelligence - Introduction to Reinforcement Learning ENSISA 2A

Ali El Hadi ISMAIL FAWAZ

November 3, 2025



Une école d'ingénieurs de l'Université de Haute-Alsace

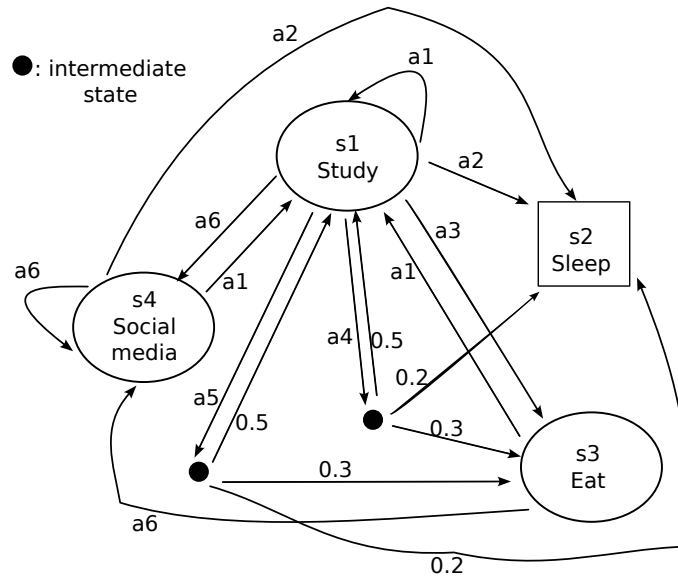


## Problem I: Covid Lockdown MDP

Consider the state-transition diagram of your Covid lockdown routine that includes 4 states:

- State  $s_1$ : Study
- State  $s_2$ : Sleep (terminal state)
- State  $s_3$ : Eat
- State  $s_4$ : Social Media

Here is a figure representing the MDP:



For this same MDP, the following actions are possible to take with their corresponding reward:

- Action  $a_1$ : Study, reward = +1
- Action  $a_2$ : Sleep, reward = +2
- Action  $a_3$ : Eat, reward = +1
- Action  $a_4$ : Sport, reward = +1
- Action  $a_5$ : Chess, reward = +1
- Action  $a_6$ : Sleep, reward = -1

Example:

If you are at state  $s_1$ , you can do the following:

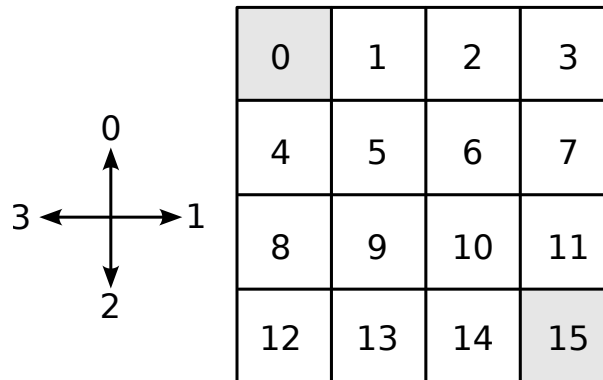
- Take action  $a_2$  and get to state  $s_2$  with a transition probability  $P_{s_1 s_2}^{a_2} = 1$  and get a reward  $R_{s_1}^{a_2} = +2$
- Take action  $a_5$ , get reward  $R_{s_1}^{a_5} = +1$  and go to:
  - state  $s_1$  again with transition probability  $P_{s_1 s_1}^{a_5} = 0.5$
  - state  $s_3$  with transition probability  $P_{s_1 s_3}^{a_5} = 0.3$
  - state  $s_2$  with transition probability  $P_{s_1 s_2}^{a_5} = 0.2$
- ...

## Questions

1. For each action, define the state transition probability matrices
2. Recall the definition of the state-value function  $V^\pi(s)$  following a policy  $\pi$  for a discounted MDP
3. Recall the Bellman expectation equation for the state-value function
4. Assuming a uniform random policy  $\pi(a|s)$  that gives an equal probability for all actions **that could be taken** from state  $s$ . For instance, in the case of state  $s_1$ , given that all possible actions can be taken, then  $\pi(a_i|s_1) = 1/6 \ i \in \{1, 2, 3, 4, 5, 6\}$ . Given the initial state-value function  $V_1^\pi(s) = 2\forall s$ , apply one epoch of the policy evaluation to compute the new state-value function  $V_2^\pi(s)\forall s$ . Assume the discount factor  $\gamma = 1$ .
5. Apply one epoch of the policy improvement algorithm and in a greedy manner, compute the new policy  $\pi_2$ .
6. Recall the definition of the optimal state-value function  $V^*(s)$  of an MDP.
7. Recall the Bellman optimality equation for state-value function.
8. Starting with initial value function  $V_1(s) = 2\forall s$ , apply one epoch using value iteration to compute  $V_2(s)\forall s$ , use discount factor  $\gamma = 1$ .
9. Is your new value function optimal ? Why ?

## Problem II: Grid World

For this exercise, you will implement the two model-based dynamic programming algorithms: value iteration and policy iteration, to solve the Grid World problem. Consider having an  $n \times n$  grid, example in Figure below when  $n = 4$ .



The position of an agent can be from squares 1 to 14 and the grey squares are the terminal state. The goal of the agent is to find the most optimal way to get to one of the two terminal states.

The actions that can be taken by the agent are:

- UP == 0
- RIGHT == 1
- DOWN == 2
- LEFT == 3

If one action throws the agent outside the grid, then it's considered to leave the agent in his current position.

A reward of  $-1$  is given to each step the agent takes until reaching a terminal state.

For this problem, consider the following parameters:

- Use  $\gamma$  discount factors of 1, 0.9 and 0.8
- The stopping criterion is defined as  $\max_s (|V_{new}(s) - V_{old}(s)|) \leq \theta$
- Use  $\theta = 10^{-4}$

## Questions

1. Implement the value iteration algorithm to solve such problem. Return the optimal policy and corresponding value function.
2. Implement the policy evaluation and policy improvement algorithms. Return the optimal policy and corresponding value function.

# Problem I:

1)

$$P = [P_{\Delta\Delta'}] \begin{bmatrix} P_{\Delta_1\Delta_1} & P_{\Delta_1\Delta_2} & P_{\Delta_1\Delta_3} & P_{\Delta_1\Delta_4} \\ P_{\Delta_2\Delta_1} & P_{\Delta_2\Delta_2} & P_{\Delta_2\Delta_3} & P_{\Delta_2\Delta_4} \\ P_{\Delta_3\Delta_1} & P_{\Delta_3\Delta_2} & P_{\Delta_3\Delta_3} & P_{\Delta_3\Delta_4} \\ P_{\Delta_4\Delta_1} & P_{\Delta_4\Delta_2} & P_{\Delta_4\Delta_3} & P_{\Delta_4\Delta_4} \end{bmatrix} \begin{array}{l} \Delta_1: \text{study} \\ \Delta_2: \text{sleep} \\ \Delta_3: \text{eat} \\ \Delta_4: \text{social media} \end{array}$$

$\begin{matrix} 1 & 2 \\ \downarrow & \downarrow \end{matrix}$   $\text{time}$

$$\underline{a_1}: \begin{array}{c} \Delta_1 \\ \Delta_2 \\ \Delta_3 \\ \Delta_4 \end{array} \begin{bmatrix} \Delta_1 & \Delta_2 & \Delta_3 & \Delta_4 \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 \end{bmatrix} \begin{array}{c} \checkmark \\ \checkmark \\ \checkmark \\ \checkmark \end{array}$$

$$\underline{a_4}: \begin{array}{c} \Delta_1 \\ \Delta_2 \\ \Delta_3 \\ \Delta_4 \end{array} \begin{bmatrix} \Delta_1 & \Delta_2 & \Delta_3 & \Delta_4 \\ 0,5 & 0,2 & 0,3 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \begin{array}{c} \checkmark \\ \checkmark \\ \checkmark \\ \checkmark \end{array}$$

$$\underline{a_2}: \begin{array}{c} \Delta_1 \\ \Delta_2 \\ \Delta_3 \\ \Delta_4 \end{array} \begin{bmatrix} \Delta_1 & \Delta_2 & \Delta_3 & \Delta_4 \\ 0 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \end{bmatrix} \begin{array}{c} \checkmark \\ \checkmark \\ \checkmark \\ \checkmark \end{array}$$

$$\underline{a_5}: \begin{array}{c} \Delta_1 \\ \Delta_2 \\ \Delta_3 \\ \Delta_4 \end{array} \begin{bmatrix} \Delta_1 & \Delta_2 & \Delta_3 & \Delta_4 \\ 0,5 & 0,2 & 0,3 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \begin{array}{c} \checkmark \\ \checkmark \\ \checkmark \\ \checkmark \end{array}$$

$$\underline{a_3}: \begin{array}{c} \Delta_1 \\ \Delta_2 \\ \Delta_3 \\ \Delta_4 \end{array} \begin{bmatrix} \Delta_1 & \Delta_2 & \Delta_3 & \Delta_4 \\ 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \begin{array}{c} \checkmark \\ \checkmark \\ \checkmark \\ \checkmark \end{array}$$

$$\underline{a_6}: \begin{array}{c} \Delta_1 \\ \Delta_2 \\ \Delta_3 \\ \Delta_4 \end{array} \begin{bmatrix} \Delta_1 & \Delta_2 & \Delta_3 & \Delta_4 \\ 0 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 1 \end{bmatrix} \begin{array}{c} \checkmark \\ \checkmark \\ \checkmark \\ \checkmark \end{array}$$

2) state-value function:  $V_{\pi}(s) = E_{\pi} [G_t | s_t = s]$

$$G_t = r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots$$

$\gamma \rightarrow$  discount factor

3) Bellman expectation equation for  $V_{\pi}(s)$ :

$$V_{\pi}(s) = E_{\pi} [r_{t+1} + \gamma V_{\pi}(s_{t+1}) | s_t = s]$$

— : variables

$$\underset{\substack{\uparrow \\ \text{policy} \\ \text{action} \rightarrow a}}{=} \sum_a \pi(a|s) \left[ \underset{\substack{\uparrow \\ \text{reward}}}{R^a} + \gamma \sum_{s'} \underset{\substack{\uparrow \\ \text{state} \\ s'}}{P_{s,s'}^a} \underset{\substack{\uparrow \\ \text{probabilité} \\ \text{de transition}}}{V_{\pi}(s')} \right]$$

Esperance  $\rightarrow$  moyenne pondérée

4)  $\pi(a|s)$  equiprobable

$s_1 \rightarrow 1/6 \quad s_2 \rightarrow 1 \quad s_3 \rightarrow 1/2 \quad s_4 \rightarrow 1/3$

$V_1^{\pi}(s) = 2 \forall s$

•  $V_2^{\pi}(s_1) = \sum_a \pi(a|s_1) [R^a + \gamma \sum_{s'} P_{s,s'}^a V_1^{\pi}(s')]$

$1/6 \rightarrow$   $\pi(a_1|s_1)$   $\downarrow$   $R^{a_1}$   $\downarrow$   $P_{s_1,s_1}^{a_1} = 1$   $\uparrow$   $V_1^{\pi}(s_1)$   $\uparrow$  initialisation pour 1<sup>ère</sup> itération

$= \pi(a_1|s_1) [R^{a_1} + V_1^{\pi}(s_1)] + P_{s_1,s_2}^{a_2} = 1$

$\pi(a_2|s_1) [R^{a_2} + V_1^{\pi}(s_2)] + \vdots$

$\pi(a_3|s_1) [R^{a_3} + V_1^{\pi}(s_3)] +$

$\pi(a_4|s_1) [R^{a_4} + P_{s_1,s_1}^{a_4} V_1^{\pi}(s_1) + P_{s_1,s_2}^{a_4} V_1^{\pi}(s_2) + P_{s_1,s_3}^{a_4} V_1^{\pi}(s_3)] +$

$\downarrow$  of matrix

$$\pi(a_5 | s_1) [R^{a_5} + P_{s_1 s_1}^{a_5} V_1^\pi(s_1) + P_{s_1 s_2}^{a_5} V_1^\pi(s_2) + P_{s_1 s_3}^{a_5} V_1^\pi(s_3)] + \pi(a_6 | s_1) [R^{a_6} + V_1^\pi(s_1)] = \frac{17}{6}$$

$$V_2(s_2) = V_1(s_2) = 2 \quad \leftarrow \text{car on ne peut pas en sortir}$$

$$V_2(s_3) = \pi(a_1 | s_3) [R^{a_1} + V_1(s_1)] + \pi(a_6 | s_3) [R^{a_6} + V_1(s_4)] = 2$$

$$V_2(s_4) = 8/3 \quad \leftarrow \text{même calcul}$$

$$5) \text{ argmax: } [1, 5, -1, 10] \quad \text{max} = 10 \\ i = 0, 1, 2, 3 \quad \text{argmax} = 3$$

$$R^{a_1} + V_2^\pi(s_1) = 1 + \frac{17}{6} = \frac{23}{6} \approx 3,83$$

$$R^{a_2} + V_2^\pi(s_2) = 2 + 2 = 4$$

$$R^{a_3} + V_2^\pi(s_3) = 1 + 2 = 3$$

$$R^{a_4} + P_{s_1 s_1}^{a_4} V_2^\pi(s_1) + P_{s_1 s_2}^{a_4} V_2^\pi(s_2) + P_{s_1 s_3}^{a_4} V_2^\pi(s_3) = 1 + 0,5 \cdot \frac{17}{6} + 0,2 \cdot 2 + 0,3 \cdot 2 = \frac{41}{12} \approx 3,42$$

$$R^{a_5} + P_{s_1 s_1}^{a_5} V_2^\pi(s_1) + P_{s_1 s_2}^{a_5} V_2^\pi(s_2) + P_{s_1 s_3}^{a_5} V_2^\pi(s_3) = 1 + 0,5 \cdot \frac{17}{6} + 0,2 \cdot 2 + 0,3 \cdot 2 = \frac{41}{12} \approx 3,42$$

$$R^{a_6} + V_2^\pi(s_4) = -1 + \frac{8}{3} = \frac{5}{3} \approx 1,67$$

$$\hookrightarrow \text{min} = 1 \Rightarrow \text{argmax}(2)$$

$$\pi_2(s_1) = a_2$$

$$R^{a_1} V_2(s_1) = 1 + 17/6 = 23/6 \approx 3,83$$

$$R^{a_2} V_2(s_1) = -1 + 8/3 = 5/3 \approx 1,67$$

$$\hookrightarrow \max = 3,83 \Rightarrow \arg\max(1)$$

$$\pi_2(s_3) = a_1$$

$$\pi_2(s_4) = a_2 \leftarrow \text{même chose}$$

$$6) V^*(s) = \max_{\pi} V^{\pi}(s)$$

$\uparrow$  maximise son policy pour  $V^{\pi}(s)$   
 state-value function optimale

$$7) V^*(s) = \max_a \left[ R^a + \gamma \sum_{s'} P_{ss'}^a V^{\pi}(s') \right]$$

$$8) \quad \left. \begin{array}{l} 1 + 2 = 3 \\ 2 + 2 = 4 \\ 1 + 2 = 3 \\ 1 + 2 = 3 \\ 1 + 2 = 3 \\ -1 + 2 = 1 \end{array} \right\}$$

$$\max = 4 = V_2^*(s_1)$$

$$V_2^*(s_2) = \cancel{\emptyset}$$

$$V_2^*(s_3) = 3$$

$$V_2^*(s_4) = 4$$

$\Rightarrow$  si on refait une autre itération et que les valeurs restent les mêmes  $\Rightarrow$  optimalité



$$g) \quad \overset{4}{V}_3^*(a_1) = 5 \neq V_3^*(a_1) = 4 \quad \Rightarrow \text{gap error optimal}$$