Introduction to Artificial Intelligence (INFO8006)

Exercises 7 – Reinforcement learning

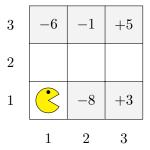
December 8, 2021

Learning outcomes

At the end of this session you should be able to

- differentiate passive and active RL;
- define and apply direct utility estimation and temporal-difference learning;
- define and apply Q-learning.

Exercise 1 Passive RL



Consider the grid-world given above and an agent who is trying to learn the optimal policy. The agent starts from the bottom-left corner and can take the actions north (N), south (S), west (W) and east (E). Rewards are only awarded for reaching the terminal (shaded) states. You observe the following trials, whose trajectories are sequences of tuples $(s_t^i, r_t^i, a_t^i, s_{t+1}^i)$.

t	Trial 1	Trial 2	Trial 3	Trial 4	Trial 5
0	(1,1),0,N,(1,2)	(1,1), 0, N, (1,2)	(1,1), 0, N, (1,2)	(1,1), 0, N, (1,2)	(1,1), 0, N, (1,2)
1	(1,2), 0, E, (2,2)	(1,2), 0, E, (2,2)	(1,2), 0, E, (2,2)	(1,2), 0, E, (2,2)	(1,2), 0, E, (2,2)
2	(2,2), 0, N, (2,3)	(2,2), 0, E, (3,2)	(2,2), 0, S, (2,1)	(2,2), 0, E, (3,2)	(2,2), 0, E, (3,2)
3	$(2,3),-1,\varnothing,\varnothing$	(3,2), 0, N, (3,3)	$(2,1), -8, \varnothing, \varnothing$	(3,2), 0, W, (2,2)	(3,2), 0, S, (3,1)
4		$(3,3),+5,\varnothing,\varnothing$		(2,2), 0, N, (2,3)	$(3,1), +3, \varnothing, \varnothing$
5				$(2,3),-1,\varnothing,\varnothing$	

Assuming a discount factor $\gamma = 1$,

1. Perform direct utility estimation of the expected utilities $V^{\pi}(s)$, given the four first trials.

In this setting, the transition model P(s'|s,a) and the reward function R(s) are unknown. We wish to learn the expected utility V^{π} of the policy π without modeling the environment, *i.e.* without building approximates $\hat{P}(s'|s,a)$ and $\hat{R}(s)$. The principle of direct utility estimation is to approximate $V^{\pi}(s)$ by the average utility $\hat{V}(s)$ of the state s within all trial trajectories, *i.e.*

$$\hat{V}(s) = \frac{1}{N(s)} \sum_{(i,j) \in I(s)} \sum_{t=0}^{\infty} \gamma^t r_{j+t}^i \approx V^{\pi}(s) = \mathbb{E}\left[\sum_{t=0}^{\infty} \gamma^t R(s_t)\right] \Big|_{s_0 = s}$$

where $I(s) = \{(i,j) : s^i_j = s\}$ indexes the occurrences of s and N(s) = |I(s)| is the number

of occurrences of s. For examples, in our case and excluding the fifth trial,

$$\begin{split} &I((1,1)) = \{(1,0),(2,0),(3,0),(4,0)\} \\ &\hat{V}((1,1)) \approx \frac{1}{4} \Big(-1\,\gamma^3 + 5\,\gamma^4 - 8\,\gamma^3 - 1\,\gamma^5 \Big) = \frac{-5}{4} \\ &I((2,2)) = \{(1,2),(2,2),(3,2),(4,2),(4,4)\} \\ &\hat{V}((2,2)) \approx \frac{1}{5} \Big(-1\,\gamma^1 + 5\,\gamma^2 - 8\,\gamma^1 - 1\,\gamma^3 - 1\,\gamma^1 \Big) = \frac{-6}{5} \\ &I((2,3)) = \{(1,3),(4,5)\} \\ &\hat{V}((2,3)) \approx \frac{1}{2} \Big(-1\,\gamma^0 - 1\,\gamma^0 \Big) = -1. \end{split}$$

Importantly, since the agent hasn't reached the states (1,3) and (3,3) yet, it is not possible to estimate their expected utility. Instead, we assume a default value of 0.

2. Update the estimated expected utilities with respect to the fifth trial using temporal-difference learning. Assume a learning rate $\alpha = 0.5$.

Direct utility estimation misses the fact that the expected utilities for different states are not independent, since they obey the Bellman equation for a fixed policy

$$V^{\pi}(s) = R(s) + \gamma \sum_{s'} P(s'|s, \pi(s)) V^{\pi}(s').$$

Then, we would like our estimates $\hat{V}(s)$ to verify the Bellman equation, *i.e.* to minimize

$$L = \underset{s'|s}{\mathbb{E}} \left[\left(R(s) + \gamma \hat{V}(s') - \hat{V}(s) \right)^2 \right].$$

One way to reach this objective is to follow the opposite of its gradient

$$\nabla L = 2 \mathop{\mathbb{E}}_{s'|s} \left[\hat{V}(s) - R(s) - \gamma \hat{V}(s') \right]$$

with respect to $\hat{V}(s)$, *i.e.* to perform *gradient descent*. Unfortunately, due to the expectation, we don't have access to the true gradient. Instead, we use a *stochastic* approximation of the gradient from observed¹ tuples (s, r, a, s') and update $\hat{V}(s)$ as

$$\hat{V}(s) \leftarrow \hat{V}(s) - \alpha \Big(\hat{V}(s) - r - \gamma \hat{V}(s')\Big)$$

$$\leftarrow (1 - \alpha)\hat{V}(s) + \alpha \Big(r + \gamma \hat{V}(s')\Big)$$

where α is the learning rate. For a fixed learning rate, the average of $\hat{V}(s)$ converge towards (oscillates around) the true expected utility $V^{\pi}(s)$. For a (slowly) decreasing learning rate, $\hat{V}(s)$ itself converges to $V^{\pi}(s)$. For historical reasons, this algorithm is called temporal-difference (TD) learning.

¹We are not limited to perform only one update for each tuple (s, r, a, s'). Indeed, as we assume the policy to be stationary, past trials are as likely to happen again as more recent ones. Therefore, reusing tuples several times (in other orders) is a good way to improve estimates while generating less trials.

Following the trajectory of the fifth trial, we perform the updates

$$V^{\pi}((1,1)) \leftarrow (1-\alpha)V^{\pi}((1,1)) + \alpha(0+\gamma V^{\pi}((1,2))) = \frac{1}{2}\frac{-5}{4} + \frac{1}{2}\frac{-5}{4} = \frac{-5}{4}$$

$$V^{\pi}((1,2)) \leftarrow (1-\alpha)V^{\pi}((1,2)) + \alpha(0+\gamma V^{\pi}((2,2))) = \frac{1}{2}\frac{-5}{4} + \frac{1}{2}\frac{-6}{5} = \frac{-49}{40}$$

$$V^{\pi}((2,2)) \leftarrow (1-\alpha)V^{\pi}((2,2)) + \alpha(0+\gamma V^{\pi}((3,2))) = \frac{1}{2}\frac{-6}{5} + \frac{1}{2}2 = \frac{+4}{10}$$

$$V^{\pi}((3,2)) \leftarrow (1-\alpha)V^{\pi}((3,2)) + \alpha(0+\gamma V^{\pi}((3,1))) = \frac{1}{2}2 + \frac{1}{2}0 = +1$$

$$V^{\pi}((3,1)) \leftarrow (1-\alpha)V^{\pi}((3,1)) + \alpha(3+0) = \frac{1}{2}0 + \frac{1}{2}3 = \frac{+3}{2}.$$

3	0	-1	+5
2	$\frac{-49}{40}$	<u>+4</u> 10	+1
1	$\frac{-5}{4}$	-8	$\frac{+3}{2}$
	1	2	3

Exercise 2 Q-learning

An agent is in an unknown environment where there are three states $\{A, B, C\}$ and two actions $\{0, 1\}$. We are given the following tuples (s, a, r, s'), generated by taking actions in the environment.

s	a	r	s'
\overline{A}	0	+2	\overline{A}
C	1	-2	A
B	1	+1	B
A	0	-1	B
B	1	-2	C
C	0	+4	B
B	0	+1	A

Assuming a discount factor $\gamma = 0.5$ and a learning rate $\alpha = 0.75$,

1. Apply the Q-learning algorithm to obtain state-action-value Q(s,a) estimates are initialized to 0.

As seen in the previous exercise, given trial trajectories of some policy π , it is possible to estimate its expected utility V^{π} even without knowing the transition model P(s'|s,a) or the reward function R(s) of the environment. However, estimating V^{π} only evaluates the quality of π , but does not describe how to improve it.

By definition, the optimal policy π^* is the one that maximizes the expect utility of the state s, i.e.

$$\pi^*(s) = \arg\max_{a} \sum_{s'} P(s'|s, a) V(s')$$

where $V(s) = \max_{\pi} V^{\pi}(s)$ is the state-value of s. Unfortunately, even if we knew V, we could not find the optimal actions without knowing the transition model. However, if we knew the state-action-value Q(s, a) of taking action a in state s, we would be able to select the optimal action

$$\pi^*(s) = \arg\max_a Q(s, a),$$

where

$$Q(s, a) = R(s) + \gamma \sum_{s'} P(s'|s, a) V(s')$$

= $R(s) + \gamma \sum_{s'} P(s'|s, a) \max_{a'} Q(s', a')$

since $V(s) = \max_a Q(s, a)$. Fortunately, state-action-values Q(s, a) can be learned in a model-free fashion using the Q-learning algorithm. Similarly to temporal-difference learning, in Q-learning, we perform stochastic gradient descent updates

$$\hat{Q}(s, a) \leftarrow \hat{Q}(s, a) - \alpha \left(\hat{Q}(s, a) - r - \gamma \max_{a'} \hat{Q}(s', a') \right)$$

$$\leftarrow (1 - \alpha) \hat{Q}(s, a) + \alpha \left(r + \gamma \max_{a'} \hat{Q}(s', a') \right)$$

from observed tuples (s, r, a, s').

In our case, the sequence of updates would be

$$\begin{split} \hat{Q}(A,0) &\leftarrow (1-\alpha)\hat{Q}(A,0) + \alpha \Big(2 + \gamma \max\{\hat{Q}(A,0),\hat{Q}(A,1)\}\Big) = \frac{1}{4}0 + \frac{3}{4}(2 + \frac{1}{2}0) = \frac{+3}{2} \\ \hat{Q}(C,1) &\leftarrow (1-\alpha)\hat{Q}(C,1) + \alpha \Big(-2 + \gamma \max\{\hat{Q}(A,0),\hat{Q}(A,1)\}\Big) = \frac{1}{4}0 + \frac{3}{4}(-2 + \frac{1}{2}\frac{3}{2}) = \frac{-15}{16} \\ \hat{Q}(B,1) &\leftarrow (1-\alpha)\hat{Q}(B,1) + \alpha \Big(2 + \gamma \max\{\hat{Q}(B,0),\hat{Q}(B,1)\}\Big) = \frac{1}{4}0 + \frac{3}{4}(1 + \frac{1}{2}0) = \frac{+3}{4} \\ \hat{Q}(A,0) &\leftarrow (1-\alpha)\hat{Q}(A,0) + \alpha \Big(-1 + \gamma \max\{\hat{Q}(B,0),\hat{Q}(B,1)\}\Big) = \frac{1}{4}\frac{3}{2} + \frac{3}{4}(-1 + \frac{1}{2}\frac{3}{4}) = \frac{-3}{32} \\ \hat{Q}(B,1) &\leftarrow \dots \end{split}$$

It should be noted that, as for TD learning, we are not limited to perform only one update for each tuple. Reusing tuples several times (in other orders) is a good way to reach convergence faster while generating less trials.

2. We now switch to a feature-based estimator $\hat{Q}(s, a) = w_0 + w_1 f_1(s, a)$, with $f_1(s, a) = 2a - 1$. Starting from weights $w_0 = w_1 = 0$, update the weights according to the approximate Q-learning algorithm.

We would like our estimates $\hat{Q}(s,a)$ to satisfy the Bellman equation, *i.e.* to minimize

$$L = \underset{s'|s,a}{\mathbb{E}} \left[\left(R(s) + \gamma \max_{a'} \hat{Q}(s', a') - \hat{Q}(s, a) \right)^2 \right].$$

To reach this objective we follow the opposite of its gradient

$$\nabla_w L = 2 \underset{s'|s,a}{\mathbb{E}} \left[\left(\hat{Q}(s,a) - R(s) - \gamma \max_{a'} \hat{Q}(s',a') \right) \nabla_w \hat{Q}(s,a) \right]$$

with respect to the weights² w of $\hat{Q}(s, a)$. Once again, due to the expectation, we are forced to use a stochastic approximation of the gradient from observed tuples (s, a, r, s') and to update the weights as

$$w \leftarrow w - \alpha \left(\hat{Q}(s, a) - r - \gamma \max_{a'} \hat{Q}(s', a')\right) \nabla_w \hat{Q}(s, a),$$

where α is the learning rate. In our case,

$$\nabla_{w}\hat{Q}(s,a) = \begin{pmatrix} \partial_{w_0} \\ \partial_{w_1} \end{pmatrix} (w_0 + w_1 f_1(s,a)) = \begin{pmatrix} 1 \\ f_1(s,a) \end{pmatrix} = \begin{pmatrix} 1 \\ 2a - 1 \end{pmatrix}$$

and the two first weight updates would be

$$\begin{split} w &\leftarrow w - \alpha \Big(\hat{Q}(A,0) - 2 - \gamma \max\{ \hat{Q}(A,0), \hat{Q}(A,1) \} \Big) \nabla_w \hat{Q}(A,0) \\ &\leftarrow \begin{pmatrix} 0 \\ 0 \end{pmatrix} - \frac{3}{4} \Big(-2 - \frac{1}{2} \max\{0,0\} \Big) \begin{pmatrix} 1 \\ -1 \end{pmatrix} = \begin{pmatrix} \frac{+3}{2} \\ \frac{-3}{2} \end{pmatrix} \\ w &\leftarrow w - \alpha \Big(\hat{Q}(C,1) + 2 - \gamma \max\{ \hat{Q}(A,0), \hat{Q}(A,1) \} \Big) \nabla_w \hat{Q}(C,1) \\ &\leftarrow \begin{pmatrix} \frac{+3}{2} \\ \frac{-3}{2} \end{pmatrix} - \frac{3}{4} \Big(2 - \frac{1}{2} \max\{0,3\} \Big) \begin{pmatrix} 1 \\ 1 \end{pmatrix} = \begin{pmatrix} \frac{+9}{8} \\ \frac{-15}{8} \end{pmatrix}. \end{split}$$

²It should be noted that the target $R(s) + \gamma \max_{a'} \hat{Q}(s', a')$ should not affected by the gradient operation as it is considered constant in the objective. In practice, this is not true as modifying the weights in a parametric function \hat{Q} does so for all state-action pairs (s, a) at once. This generally makes the approximate Q-learning algorithm very unstable.

Exercise 3 Football

ULiège's football team is playing against UCL's team next week. With a lot of losses this season, Liège needs to improve their attack strategy to win the game and increase their popularity. Luckily, the team captain follows INFO8006 and knows how to model the attack as a Markov Decision Process. The captain considers four states close, away, fail, and goal and two actions pass and shoot. Although the transition probabilities are unsure, the possible transitions (s, a, s') are known. To each transition is associated an increase/decrease of the team's popularity.

s	a	s'	R(s,a,s')
close	pass	close	+1
close	pass	away	-1
close	pass	fail	-2
close	shoot	close	+3
close	shoot	fail	-5
close	shoot	goal	+10
away	pass	close	+2
away	pass	away	0
away	pass	fail	-3
away	shoot	close	+3
away	shoot	fail	-10
away	shoot	goal	+20

The current strategy of the team is to always shoot. Last match, they had several attack opportunities, resulting in the following actions.

s	a	s'
close	shoot	goal
close	shoot	close
close	shoot	goal
close	shoot	fail
away	shoot	fail
away	shoot	close
away	shoot	fail
away	shoot	fail

Assuming a discount factor $\gamma = 0.75$ and a learning rate $\alpha = 0.25$,

- 1. Build an estimator $\hat{P}(s'|s,a)$ of the transition model and, from it, determine the expected utility V^{π} of the team's current policy π .
- 2. Perform direct utility estimation of the expected utility V^{π} . Do you observe a difference with the previous estimation? Why?

The captain found the tapes of the previous season where they had much more success. Together with the team, the captain selects the following instructive actions.

s	a	s'
close	pass	fail
close	pass	close
close	pass	away
close	pass	fail
away	pass	away
away	pass	close
away	pass	close
away	pass	fail

- 3. Apply the Q-learning algorithm to obtain state-action-value Q(s,a) estimates. Estimates are initialized to 0.
- 4. Determine the optimal policy according to the state-action-value estimates.

Supplementary materials

• Playing Atari with Deep Reinforcement Learning



• Chapter 21 of the reference textbook.