

# Artificial Neural Networks (Gerstner). Exercises for week 4

## Variants of TD-learning methods and eligibility traces

### Exercise 0. Recap: SARSA on Linear Track<sup>1</sup>

Please make sure that you have finished ‘Exercise 2. SARSA algorithm’ from week 2 before you continue. Exercise 2 from week 2 is important for the lecture today at 14h15. The learning outcomes of the SARSA exercise are (i) to be at ease using the SARSA update formula and (ii) to understand why learning with SARSA takes many episodes.

### Exercise 1. Expected SARSA and a new variant of TD learning<sup>2</sup>

We have seen the algorithm ‘Expected Sarsa’ and ‘TD-algorithm in the narrow sense’. Suppose you invent a new algorithm that keeps tables of both  $Q$ -values and  $V$ -values. To do so consider the following:

- Rewrite the updated step of ‘Expected Sarsa’ using both  $Q$ -values and  $V$ -values in the update rule.
- Write down the full pseudo-algorithm.

*Hint:* (i) keep track of the sequence of update steps of  $Q$ -values and  $V$ -values in your new algorithm. (ii) At which point in the sequence of online steps through the graph can you update  $Q(s, a)$  and how far do you have to look backward?

- Sketch the ‘back-up-diagram’ of your algorithm and compare it to that of SARSA. What are the costs and benefits of your new algorithm compared to SARSA?

### Exercise 2. Monte Carlo versus expected SARSA<sup>3</sup>

We compare two algorithms: The Batch-Monte-Carlo algorithm (around slide 47), and the Online-Expected-SARSA (around slide 19). The aim of the exercise is to understand why bootstrap algorithms (i.e., TD algorithms) perform, under some conditions, more efficiently than a naive estimation of  $Q$ -Values via Monte-Carlo algorithms. We set the discount factor to  $\gamma = 1$  and run 5 episodes in a given environment:

#### 5 episodes, first action is always $a_1$ .

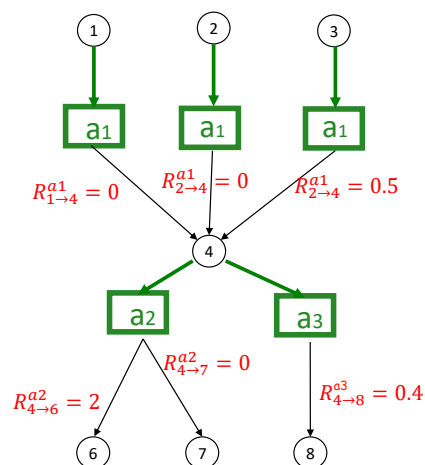
Episode 1: States 1-4-7 with action  $a_2$ , Return=0

Episode 2: States 1-4-8 with action  $a_3$ , Return=0.4

Episode 3: States 2-4-6 with action  $a_2$ , Return=2

Episode 4: States 2-4-8 with action  $a_3$ , Return=0.4

Episode 5: States 3-4-7 with action  $a_2$ , Return=0.5



Each episode starts in one of the states 1 or 2 or 3 with action  $a_1$ . In state 4 there is a choice between actions  $a_2$  and  $a_3$  which are taken with equal probability  $\pi = 0.5$ . The transition sequence and total

<sup>1</sup>The result of Exercise 0 will be used in the second lecture of this week.

<sup>2</sup>The result of Exercise 1 will be used in the second lecture of this week.

<sup>3</sup>The result of Exercise 2 will be used in the second lecture of this week.

return for each of the 5 episodes is given in the figure above. All rewards are deterministic and only depend on the transition  $(s, a, s')$ .

- Calculate the Q-values in state 1, 2, 3, and 4 using Batch-Monte-Carlo (i.e., average total returns from each starting state).
- Calculate the Q-values in state 1, 2, 3, and 4 using Online-Expected SARSA. For a given Q-value  $Q(s, a)$ , use  $\eta_1 = 1$  the FIRST TIME you update this value and  $\eta \in [0, 0.5]$  for all later update steps.

*Hint:* You can neglect terms of order  $\eta^2$ .

- You can choose the initial state for episode 6. Which initial state looks best in case a (Batch-Monte-Carlo)? Which initial state looks best in case b (Online-Expected SARSA)?

Where does the difference come from? Which solution is closer to the exact Q-values under the  $\pi = 0.5$  policy in state 4? Why?

- What happens qualitatively if the order of episodes 1 and 2 were switched? What would be the value of  $Q(s = 1, a_1)$  in this case? How would the other Q-values change?

### Exercise 3. Eligibility traces

In week 2 in exercise 2, we applied the SARSA algorithm to the case of a linear track with actions ‘up’ and ‘down’. We found that it takes a long time to propagate the reward information through state space. The eligibility trace is introduced as a solution to this problem.

Reconsider the linear maze from Figure 1 in exercise 2, but include an eligibility trace: for each state  $s$  and action  $a$ , a memory  $e(s, a)$  is stored. At each time step, all the memories are reduced by a factor  $\lambda < 1$ :  $e(s, a) = \lambda e(s, a)$ , except for the memory corresponding to the current state  $s^*$  and action  $a^*$ , which is incremented:

$$e(s^*, a^*) = \lambda e(s^*, a^*) + 1. \quad (1)$$

Now, unlike the case without eligibility trace, all Q-values are updated at each time step according to the rule

$$\forall (s, a) \quad \Delta Q(s, a) = \eta [r - (Q(s^*, a^*) - Q(s', a'))] e(s, a). \quad (2)$$

where  $s^*, a^*$  are the current state and action, and  $s', a'$  are the immediately following state and action.

We want to check whether the information about the reward propagates more rapidly. To find out, assume that the agent goes straight down in the first trial. In the second trial it uses a greedy policy. Calculate the Q-values after two complete trials and report the result.

Hint: Reset the eligibility trace to zero at the beginning of each trial.

### Exercise 4. Eligibility traces in continuous space

The left part of Figure 1 shows a different representation of last week’s “linear track” exercise: the vertical divisions represent different states, and the two columns correspond to the two possible actions available to the agent: go up or down. Each square represents a possible state-action combination, and thus a Q-value. (Note that the uppermost “up” action and the lowermost “down” action should be “greyed out”: they are impossible. But this is not relevant to the rest of this exercise.)

Suppose now that the agent moves in a continuous 1-dimensional space  $0 \leq x \leq 1$ , with the target located at  $x = 0$ . Separate this state space into  $n$  equal bins of width  $\Delta x = 1/n$ . In each time step, the agent moves by one bin. Vary the discretization by varying  $n$ :  $n = 4, 8, 16 \dots$

- Suppose that one action (such as move down) corresponds to one time step  $\Delta t$  in ‘real time’. How should we rescale the parameter  $\Delta t$ , so that the speed  $v = \Delta x / \Delta t$  remains constant when we change the discretization?

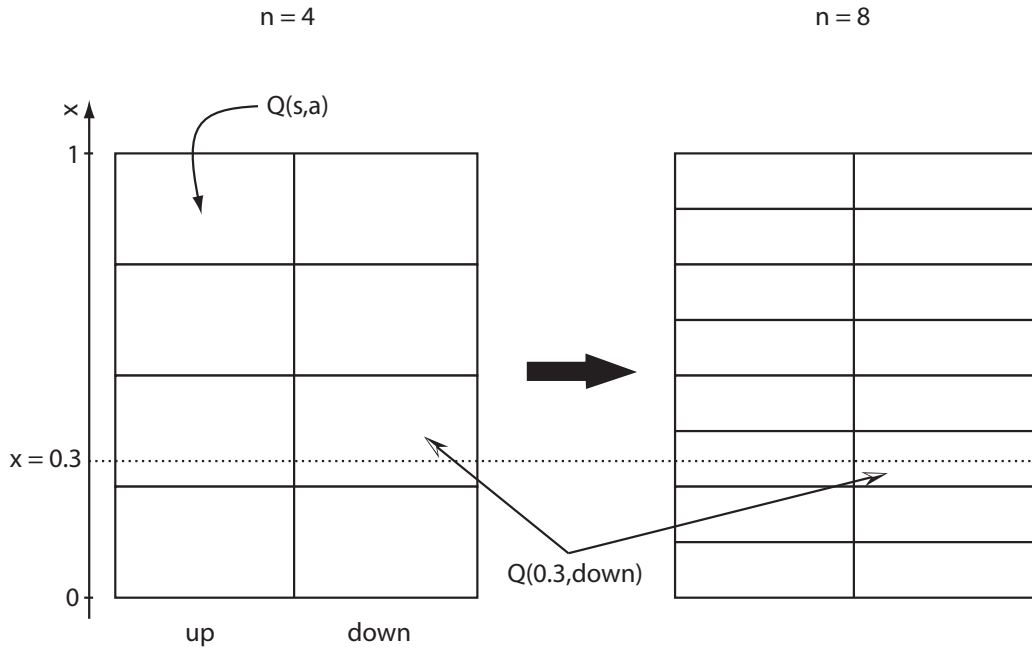


Figure 1: Figure for Exercise 4

- b. We use an eligibility trace with decay parameter  $\lambda$ . How should we rescale  $\lambda$ , in order that the “speed of information propagation” in SARSA( $\lambda$ ) remains constant?

Hint: Consider the Q-value at a fixed  $x$ , for example at  $x = 0.5$ , after 2 complete learning trials.

### Exercise 5. 3-step SARSA algorithm

In class we have discussed the SARSA algorithm and shown that, after convergence, the resulting Q-values solve (in expectation) the Bellman equation for *neighboring* states (Variant A in the slides, fixed/non-fluctuating Q-values after convergence). Your friend claims that a 3-step SARSA for

$$\Delta Q(s_t, a_t) = \eta [r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \gamma^3 Q(s_{t+3}, a_{t+3}) - Q(s_t, a_t)] , \quad (3)$$

should work just as well.

To simplify the analysis, we assume that the environment has no loops (i.e., the graph is directed) so that we can consider  $\gamma = 1$ .

- a. Assume that the 3-step SARSA algorithm converges in expectation. Proceed as during the lecture to show that  $\mathbb{E}[\Delta Q(s_t, a_t)] = 0$  implies

$$Q(s_t, a_t) = \sum_{s'} P_{s_t \rightarrow s'}^{a_t} \left[ R_{s_t \rightarrow s'}^{a_t} + \sum_{a'} \pi(s', a') B_1(s', a') \right]$$

where

$$B_1(s', a') = \sum_{s''} P_{s' \rightarrow s''}^{a'} \left[ R_{s' \rightarrow s''}^{a'} + \sum_{a''} \pi(s'', a'') B_2(s'', a'') \right]$$

$$B_2(s'', a'') = \sum_{s'''} P_{s'' \rightarrow s'''}^{a''} \left[ R_{s'' \rightarrow s'''}^{a''} + \sum_{a'''} \pi(s''', a''') Q(s''', a''') \right]$$

- b. Show the equivalence of the previous equation to the 1-step Bellman equation.

**Exercise 6. Computer exercises: Environment 1 (part 2)**<sup>1</sup>

Complete the computer exercise for environment 1.

---

<sup>1</sup>Start this exercise in the second exercise session of this week.