

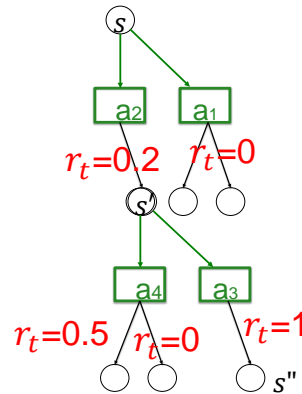
# Artificial Neural Networks (Gerstner). Solutions for week 3

## Variations of SARSA and continuous space

### Exercise 1. Monte-Carlo Batch versus Bootstrap Batch

example trials:

- 1:  $\mathbf{s}, \mathbf{a}_2 \rightarrow r=0.2, \mathbf{s}', \mathbf{a}_4 \rightarrow r=0$
- 2:  $\mathbf{s}', \mathbf{a}_3 \rightarrow r=1$
- 3:  $\mathbf{s}', \mathbf{a}_4 \rightarrow r=0$
- 4:  $\mathbf{s}', \mathbf{a}_3 \rightarrow r=1$
- 5:  $\mathbf{s}, \mathbf{a}_1 \rightarrow r=0$
- 6:  $\mathbf{s}', \mathbf{a}_4 \rightarrow r=0$
- 7:  $\mathbf{s}', \mathbf{a}_4 \rightarrow r=0.5$
- 8:  $\mathbf{s}', \mathbf{a}_3 \rightarrow r=1$
- 9:  $\mathbf{s}, \mathbf{a}_2 \rightarrow r=0.2, \mathbf{s}', \mathbf{a}_4 \rightarrow r=0.5$
- 10:  $\mathbf{s}, \mathbf{a}_1 \rightarrow r=0$



The aim is to estimate Q-values on the graph as shown in the figure on top. We run 10 exploration trials with the rewards as indicated in the figure. Some trials start in state  $s$  others in state  $s'$ . Discount factor  $\gamma$  is equal to 1, and the policy is the random policy:  $\pi(s, a) = 0.5$  for all state-action pairs. You use three different methods for estimating the Q-values: a. Batch Monte-Carlo, b. Batch-Expected SARSA, and c. Expected-SARSA.

- a. Use Monte-Carlo in batch version to calculate the four Q-value (all combination of state-action pairs). To do so, evaluate for all state-action pairs the cumulative expected reward  $R(s, a)$  with discount factor  $\gamma = 1$  by averaging over all trials that either start at  $(s, a)$  or pass through  $(s, a)$ .
- b. Use Batch-Expected SARSA (= dynamic programming with an empirical component) to evaluate the estimated Q-values  $Q(s, a)$  by exploiting the Bellman equation with discount factor  $\gamma = 1$ .

Hints:

(i) Start at the terminal state and work up backwards.

(ii) The empirical Bellman equation is:

$$Q(s, a) = \langle r_t \rangle + \gamma \sum_{s'} P_{s \rightarrow s'}^{\text{emp}, a} \sum_{a'} \pi(s', a') Q(s', a')$$

Here  $\langle r_t \rangle$  is the empirical average over the immediate rewards for the combination  $(s, a)$  (over all trials), and  $P_{s \rightarrow s'}^{\text{emp}, a} = n(s'|s, a) / \sum_{s''} n(s''|s, a)$  is the empirical average of the transition matrix where  $n(s''|s, a)$  is the number of trials in which you have taken action  $a$  in state  $s$  and landed in  $s''$ .

(iii) Note that for action  $a_2$ , the branching ratio  $P_{s \rightarrow s'}^{\text{emp}, a}$  is unity, since there is only one possible landing state.

- c. Use Expected-SARSA in the standard *online* version to evaluate the estimated Q-values  $Q(s, a)$ . To do so, use the optimal time-dependent learning rate  $\eta$  that we have seen in exercise 1 of last week. Since we have several states, you have to make  $\eta(s, a)$  to depend

on count  $n(s, a)$  for each state-action pair  $\eta(s, a) = 1/n(s, a)$ . You initialize all Q-values at zero and use discount factor  $\gamma = 1$ .

Hints:

- (i) Convince yourself that after the first trial  $Q(s, a_2)$  increases by 0.2.
  - (ii) Use the result of last week and exploit that the results of the specific time-dependent  $\eta(s, a) = 1/n(s, a)$  in multi-arm bandit (and one-step horizon) tasks are in fact *equivalent* to batch. Thus, we know that after 8 trials, the Q-values  $Q_8(s', a_3)$  and  $Q_8(s', a_4)$  resulting from online Expected-SARSA gives are equivalent of a batch evaluation of the rewards after 8 trials.
  - (iii) Then analyze the update of  $\Delta Q(s, a_2)$  after episode number 9. Remember that this is the *second* time that you play this action.
- d. Calculate the exact Q-values for the random policy and assuming a balanced branching ratio where applicable.
- e. Which of the algorithms performs best on this example?

**Solution:**

a. For each pair of state and action, we average discounted return over all trials:

- $Q(s', a_4) = \frac{1}{5}(0 + 0 + 0 + 0.5 + 0.5) = 0.2$  (corresponding to trials 1, 3, 6, 7, and 9)
- $Q(s', a_3) = \frac{1}{3}(1 + 1 + 1) = 1$  (corresponding to trials 2, 4, and 8)
- $Q(s, a_2) = \frac{1}{2}((0.2 + 0) + (0.2 + 0.5)) = 0.45$  (corresponding to trials 1 and 9)
- $Q(s, a_1) = \frac{1}{2}(0 + 0) = 0$  (corresponding to trials 5 and 10)

b. As a result of applying Batch-Expected SARSA, the estimated Q-values are

- $Q(s', a_4) = \langle r_t \rangle = \frac{1}{5}(0 + 0 + 0 + 0.5 + 0.5) = 0.2$
- $Q(s', a_3) = \langle r_t \rangle = \frac{1}{3}(1 + 1 + 1) = 1$
- $Q(s, a_2) = \langle r_t \rangle + \frac{Q(s', a_3) + Q(s', a_4)}{2} = \frac{1}{2}(0.2 + 0.2) + \frac{1 + 0.2}{2} = 0.8$
- $Q(s', a_1) = \langle r_t \rangle = \frac{1}{2}(0 + 0) = 0$

c. For Expected-SARSA we have:

- After trial 1:  $Q(s, a_2) \leftarrow Q(s, a_2) + \frac{1}{1}RPE = 0 + 1(0.2 + 0 - 0) = 0.2$
- After trial 8:  $Q(s', a_3) = \frac{1}{3}(1 + 1 + 1) = 1$  and  $Q(s', a_4) = \frac{1}{4}(0 + 0 + 0 + 0.5) = 0.125$
- After trial 9:  $Q(s, a_2) \leftarrow Q(s, a_2) + \frac{1}{2}RPE = 0.2 + \frac{1}{2}(0.2 + \frac{0.125 + 1}{2} - 0.2) = 0.48125$  and  $Q(s', a_4) = 0.2$

Therefore, at the end we have

- $Q(s', a_4) = 0.2$
- $Q(s', a_3) = 1$

- $Q(s, a2) = 0.48125$
- $Q(s', a1) = 0$

d. The exact Q-values as the solutions to the Bellman equations corresponding to the figure are

- $Q(s', a4) = \frac{1}{2}(0 + 0.5) = 0.25$
- $Q(s', a3) = \frac{1}{1}(1) = 1$
- $Q(s, a2) = \frac{1}{1}(0.2) + \frac{1+0.25}{2} = 0.825$
- $Q(s', a1) = \frac{1}{2}(0 + 0) = 0$

e. Batch-Expected SARSA > Expected SARSA > Batch Monte-Carlo

### Exercise 2. Eligibility traces

In week 2 in exercise 3, 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 Fig 2. in exercise 3, 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  $\lambda$ :  $e(s, a) = \lambda e(s, a)$ , 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.

### Solution:

The table below shows the evolution of the Q values for each relevant state action pair during the first 2 trials, starting at the first step when there is a non-zero update  $\Delta = \eta [r - (Q(s^*, a^*) - Q(s', a'))]$ . We assume that the agent goes straight down in the first trial, that it always picks the best action, and that the eligibility traces are reset when the agent picks the reward, and the agent is put back to the starting position.

trial	transition		$(s, a_1)$	$(s', a_1)$	$(s'', a_1)$	$\Delta$
1	$s'' \rightarrow s'''$	$Q$ $e$	0 $\lambda^2$	0 $\lambda$	0 1	$\eta$ -
2	$s \rightarrow s'$	$Q$ $e$	$\eta\lambda^2$ 1	$\eta\lambda$ 0	$\eta$ 0	$\eta^2(\lambda - \lambda^2)$ -
2	$s' \rightarrow s''$	$Q$ $e$	$\eta\lambda^2 + \eta^2(\lambda - \lambda^2)$ $\lambda$	$\eta\lambda$ 1	$\eta$ 0	$\eta^2(1 - \lambda)$ -
2	$s'' \rightarrow s'''$	$Q$ $e$	$\eta\lambda^2 + 2\eta^2\lambda - 2\eta^2\lambda^2$ $\lambda^2$	$\eta\lambda + \eta^2(1 - \lambda)$ $\lambda$	$\eta$ 1	$\eta - \eta^2$ -
3	$s \rightarrow s'$	$Q$ $e$	$2\eta\lambda^2 + 2\eta^2\lambda - 3\eta^2\lambda^2$ 1	$2\eta\lambda + \eta^2 - 2\eta^2\lambda$ 0	$2\eta - \eta^2$ 0	$\dots$ -

...

Although the  $Q$ -values for  $s''$  are the same as without the eligibility trace (see exercise from last week), the  $Q$ -values for  $s$  and  $s'$  already start to approach their asymptotical value (i. e. 1) in the first trial.

### Exercise 3. Q-values for continuous states

We approximate the state-action value function  $Q(s, a)$  by a weighted sum of basis functions (BF):

$$Q(s, \tilde{a}) = \sum_j w_{\tilde{a}j} \Phi(s - s_j),$$

where  $\Phi(x)$  is the BF “shape”, and the  $s_j$ 's represent the centers of the BFs.

Calculate

$$\frac{dQ(s, a)}{dw_{ai}},$$

the gradient of  $Q(s, a)$  along  $w_{ai}$  for a specific weight linking the basis function  $i$  to the action  $a$ .

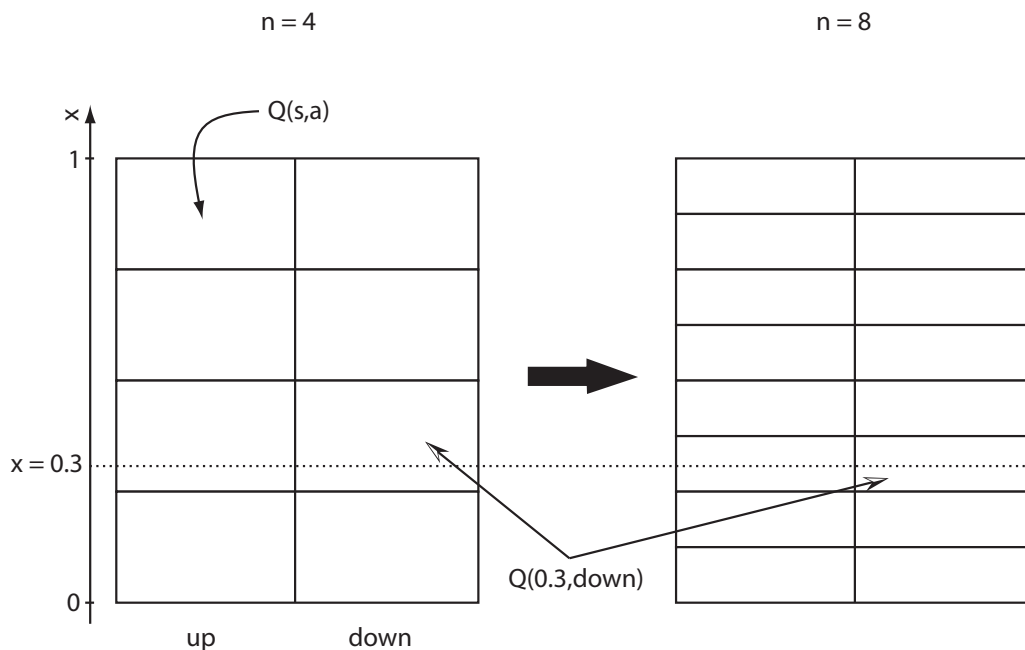
#### Solution:

Using the definition of  $Q(s, a)$  given, we find the gradient:

$$\frac{dQ(s, a)}{dw_{aj}} = \Phi(s - s_j).$$

Therefore the direction of the gradient vector ( $dQ/dw_{aj}$ ) for  $j = 1, \dots, K$  is given by the magnitude of responses ( $\Phi(s - s_j)$ ) of all basis functions.

### Exercise 4. Eligibility traces in continuous space



The left part of the figure above shows a different representation of last week’s “linear track” exercise: the vertical divisions represent different states, and the two column 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?
- 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 for example, the  $Q$ -value at  $x = 0.5$  after 2 complete learning trials.

**Solution:**

- If  $\Delta x / \Delta t$  has to remain constant, then  $\Delta t$  should vary like  $\Delta x$ , i.e.,  $\Delta t \propto 1/n$ .
- A point “sitting” at  $x$  is  $x / \Delta x = x \cdot n$  steps away from the reward (assuming we always choose the “down” action). As we have seen in exercise 2, on the first trial, the  $Q$ -value of a state  $d$  steps from the trial is updated proportional to  $\lambda^d$ . Thus, if we want the update to stay constant under rescaling, we need

$$\Delta Q(s, a) \propto \lambda^d = \lambda^{x \cdot n} = cst$$

This holds if we replace  $\lambda$  by a “rescaled”  $\tilde{\lambda} = \lambda^{\frac{1}{n}}$ .

### Exercise 5. Gradient-based learning of Q-values

Assume again that the Q-values are expressed as a weighted sum of 400 basis functions:  $Q(s, a) = \sum_{k=1}^{400} w_a^k \Phi(s - s_k)$ . For the moment the function  $\Phi$  is arbitrary, but you may think of it as a Gaussian function. Note that  $s$  and  $s_k$  are usually vectors in  $\mathbb{R}^N$  in this case. There are 3 different actions so that the total number of weights is 1200. Now consider the error function  $E_t = \frac{1}{2} \delta_t^2$ , where

$$\delta_t = r_t + \gamma \cdot Q(s', a') - Q(s, a) \quad (3)$$

is the reward prediction error. Our aim is to optimize  $Q(s, a)$  by changing the parameters  $w$ .

- a. Find a learning rule that minimizes the error function  $E_t$  by gradient decent. Consider the case where the actions  $a$  and  $a'$  are different.

Write down the learning rule. How many weights need to be updated in each time step?

- b. Find a learning rule that minimizes the error function  $E_t$  by gradient decent. Consider the case where the actions  $a$  and  $a'$  are the same.

Write down the learning rule.

Is there any difference to the case considered in a)?

- c. Suppose that the input space is two-dimensional and you discretize the input in 400 small square 'boxes' (i.e.,  $20 \times 20$ ). The basis function  $\Phi(s - s_k)$  is now the indicator function: it has a value equal to one if the current state  $s$  is in 'box'  $k$  and zero otherwise.

How do your results from (a) and (b) look like in this case?

- d. The learning rule in c) is very similar to standard SARSA. What is the difference?

- e. Assume that  $Q(s', a')$  in Equation 3 does not depend on the weights. How does the learning rule look like for a, b and c. How is your result related to standard discrete SARSA?

#### Solution:

- a. Let's start by computing the derivative of  $Q(s, a)$  with respect to  $w_k^{\tilde{a}}$  (we'll use this later):

$$\frac{dQ(s, a)}{dw_k^{\tilde{a}}} = \delta_{a\tilde{a}} \Phi(s - s_{\tilde{k}}),$$

where  $\delta_{a\tilde{a}}$  is the Kronecker, i.e., it is 1 if  $a = \tilde{a}$ , and 0 otherwise (not to be confused with  $\delta_t$ ).

We then compute the gradient, i.e., the derivative of  $E_t$  with respect to  $w_k^{\tilde{a}}$ , using the chain rule a few times and the result above:

$$\begin{aligned} \frac{dE_t}{dw_k^{\tilde{a}}} &= \delta_t \left[ \gamma \frac{dQ(s', a')}{dw_k^{\tilde{a}}} - \frac{dQ(s, a)}{dw_k^{\tilde{a}}} \right] \\ &= \delta_t [\gamma \delta_{a'\tilde{a}} \Phi(s' - s_{\tilde{k}}) - \delta_{a\tilde{a}} \Phi(s - s_{\tilde{k}})]. \end{aligned}$$

To turn this into a learning rule, we have to move the weights in the direction that *minimizes* the error, i.e.

$$\Delta w_k^{\tilde{a}} = -\eta \frac{dE_t}{dw_k^{\tilde{a}}} = \eta \delta_t [\delta_{a\tilde{a}} \Phi(s - s_{\tilde{k}}) - \gamma \delta_{a'\tilde{a}} \Phi(s' - s_{\tilde{k}})]. \quad (4)$$

$2 \cdot 400$  weights (for actions  $a$  and  $a'$ ) need to be updated in each step.

b. In the case where  $a = a'$  (i.e., the action taken is the same in the two consecutive steps):

$$\Delta w_{\tilde{k}}^{\tilde{a}} = \eta \delta_t (\Phi(s - s_{\tilde{k}}) - \gamma \Phi(s' - s_{\tilde{k}})) \delta_{a\tilde{a}}. \quad (5)$$

400 weights need to be updated.

- c. The learning rules look the same. Only two weights need to be updated if  $a \neq a'$  since the basis function is zero  $\Phi(s - s_k) = 0$  if  $s$  does not fall into the  $k$ th box. If  $a = a'$ , we have to update either one or two weights, depending on whether  $s$  and  $s'$  fall into the same box or not.
- d. Without the term  $\gamma \Phi(s' - s_{\tilde{k}})$ , the algorithm becomes identical to standard SARSA.
- e. If  $Q(s', a')$  is a fixed target that does not depend on the weights, the term  $\gamma \Phi(s' - s_{\tilde{k}})$  drops out of the learning rules and with the small square box function approximator the algorithm becomes equivalent to SARSA. This is sometimes called semigradient learning.

### Exercise 6. Consistency condition for 3-step SARSA

In class we have seen the arguments leading to the error function arising from the consistency condition of Q-values.

$$E = 0.5 \sum [\delta_t]^2$$

with  $\delta_t = r_t + \gamma Q(s', a') - Q(s, a)$  This specific consistency condition corresponds to 1-step SARSA.

Write down an analogous consistency condition for 3-step SARSA.

#### Solution:

The error function is  $E = 0.5 \sum [\delta_t]^2$ , but with a consistency condition  $\delta_t = r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \gamma^3 Q(s_{t+3}, a_{t+3}) - Q(s, a)$  where  $s_{t+3}, a_{t+3}$  are state and action three time steps after taking action  $a$  in state  $s$ . Explanation: the value  $Q(s, a)$  must be explained by the (average of the) rewards in the next three steps plus the  $Q$ -value of the action three time steps after taking action  $a$  in state  $s$ . The averaging is taken implicitly by an online on-policy algorithm once the same sequence has been taken multiple times.