

# Artificial Neural Networks (Gerstner). Solutions for week 13

## Reinforcement Learning and the Brain

### Exercise 1. A biological interpretation of the Advantage Actor-Critic with Eligibility traces

In this exercise you will show how applying Advantage Actor-Critic with eligibility traces to a softmax policy in combination with a linear read-out function leads to a biologically plausible learning rule.

Consider a policy and a value network as in Figure 1 with  $K$  input neurons  $\{y_k = f(x - x_k)\}_{k=1}^K$ . The policy network is parameterized by  $\theta$  and has three output neurons corresponding to actions  $a_1$ ,  $a_2$  and  $a_3$  with 1-hot coding. If  $a_k = 1$ , action  $a_k$  is taken. The output neurons are sampled from a softmax policy: The probability of taking action  $a_i$  is given by

$$\pi_\theta(a_i = 1|x) = \frac{\exp[\sum_k \theta_{ik} y_k]}{\sum_j \exp[\sum_k \theta_{jk} y_k]}. \quad (1)$$

In addition, consider the exponential value network  $\hat{v}_w(x) = \exp[\sum_k w_k y_k]$ .

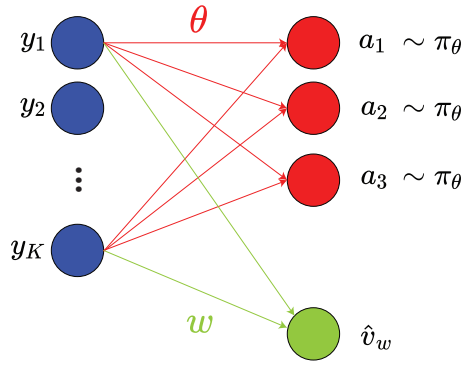


Figure 1: The network structure.

Assume the transition to state  $x^{t+1}$  with a reward of  $r^{t+1}$  after taking action  $a^t$  at state  $x^t$ . The learning rule for the Advantage Actor-Critic with Eligibility traces is

$$\begin{aligned} \delta &\leftarrow r^{t+1} + \gamma \hat{v}_w(x^{t+1}) - \hat{v}_w(x^t) \\ z^w &\leftarrow \lambda^w z^w + \nabla_w \hat{v}_w(x^t) \\ z^\theta &\leftarrow \lambda^\theta z^\theta + \nabla_\theta \pi_\theta(a^t|x^t) \\ w &\leftarrow w + \alpha^w z^w \delta \\ \theta &\leftarrow \theta + \alpha^\theta z^\theta \delta \end{aligned}$$

Your goal is to show that this learning rule applied to the network of Figure 1 has a biological interpretation.

- a. Show that

$$\frac{d}{dw_5} \hat{v}_w(x^t) = y_5^t \hat{v}_w(x^t). \quad (2)$$

- b. Interpret the update of the eligibility trace  $z_5^w$  in terms of a ‘presynaptic factor’ and a ‘postsynaptic factor’. Can the rule be implemented in biology?

- c. Show that

$$\frac{d}{d\theta_{35}} \ln[\pi_\theta(a^t|x^t)] = [a_3^t - \pi_\theta(a_3 = 1|x^t)] y_5^t. \quad (3)$$

Hint: simply insert the softmax and then take the derivative.

- d. Interpret the update of the eligibility trace  $z_{35}^\theta$  in terms of a ‘presynaptic factor’ and a ‘postsynaptic factor’. Can the rule be implemented in biology?
- e. Interpret the update of the weights  $w_5$  and  $\theta_{35}$  in the framework of three factor learning rules. Can the rule be implemented in biology?