

# COMP9444 Neural Networks and Deep Learning

## Session 2, 2018

### Solutions to Exercise 6: Reinforcement Learning

This page was last updated: 09/11/2018 08:45:06

Consider an environment with two states  $S = \{S_1, S_2\}$  and two actions  $A = \{a_1, a_2\}$ , where the (deterministic) transitions  $\delta$  and reward  $R$  for each state and action are as follows:

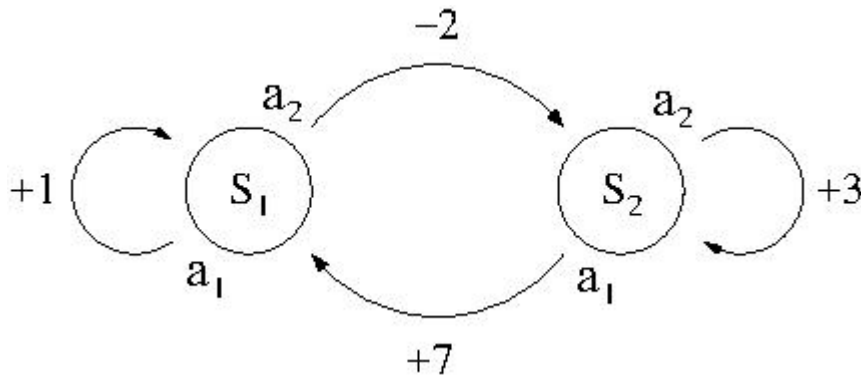
$$\delta(S_1, a_1) = S_1, R(S_1, a_1) = +1$$

$$\delta(S_1, a_2) = S_2, R(S_1, a_2) = -2$$

$$\delta(S_2, a_1) = S_1, R(S_2, a_1) = +7$$

$$\delta(S_2, a_2) = S_2, R(S_2, a_2) = +3$$

1. Draw a picture of this environment, using circles for the states and arrows for the transitions.



2. Assuming a discount factor of  $\gamma = 0.7$ , determine:

- a. the optimal policy  $\pi^* : S \rightarrow A$

$$\pi^*(S_1) = a_2$$

$$\pi^*(S_2) = a_1$$

- b. the value function  $V : S \rightarrow \mathbb{R}$

$$V(S_1) = -2 + \gamma V(S_2)$$

$$V(S_2) = +7 + \gamma V(S_1)$$

$$\text{So } V(S_1) = -2 + 7\gamma + \gamma^2 V(S_1)$$

$$\text{i.e. } V(S_1) = (-2 + 7\gamma)/(1 - \gamma^2) = (-2 + 7 \times 0.7)/(1 - 0.49) = 5.69$$

$$\text{Then } V(S_2) = 7 + 0.7 \times 5.69 = 10.98$$

- c. the "Q" function  $Q : S \times A \rightarrow \mathbb{R}$

$$Q(S_1, a_1) = 1 + \gamma V(S_1) = 4.98$$

$$Q(S_1, a_2) = V(S_1) = 5.69$$

$$Q(S_2, a_1) = V(S_2) = 10.98$$

$$Q(S_2, a_2) = 3 + \gamma V(S_2) = 10.69$$

Writing the Q values in a matrix, we have:

Q	a <sub>1</sub>	a <sub>2</sub>
S <sub>1</sub>	4.98	5.69
S <sub>2</sub>	10.98	10.69

Trace through the first few steps of the Q-learning algorithm, with a learning rate of 1 and with all Q values initially set to zero. Explain why it is necessary to force exploration through probabilistic choice of actions, in order to ensure convergence to the true Q values.

With a deterministic environment and a learning rate of 1, the Q-Learning update rule is

$$Q(S, a) \leftarrow r(S, a) + \gamma \max_b Q(\delta(S, a), b)$$

Let's assume the agent starts in state S<sub>1</sub>. Since the initial Q values are all zero, the first action must be chosen randomly. If action a<sub>1</sub> is chosen, the agent will get a reward of +1 and update

$$Q(S_1, a_1) \leftarrow 1 + \gamma \times 0 = 1$$

If we do not force exploration, the agent will always prefer action a<sub>1</sub> in state S<sub>1</sub>, and will never explore action a<sub>2</sub>. This means that Q(S<sub>1</sub>, a<sub>2</sub>) will remain zero forever, instead of converging to the true value of 5.69. If we do force exploration, the next steps may look like this:

current state	chosen action	new Q value
S <sub>1</sub>	a <sub>2</sub>	-2 + $\gamma \cdot 0 = -2$
S <sub>2</sub>	a <sub>2</sub>	+3 + $\gamma \cdot 0 = +3$

At this point, the table looks like this:

Q	a <sub>1</sub>	a <sub>2</sub>
S <sub>1</sub>	1	-2
S <sub>2</sub>	0	3

Again, we need to force exploration, in order to get the agent to choose a<sub>1</sub> from S<sub>2</sub>, and to again choose a<sub>2</sub> from S<sub>1</sub>

current state	chosen action	new Q value
S <sub>2</sub>	a <sub>1</sub>	+7 + $\gamma \cdot 1 = 7.7$
S <sub>1</sub>	a <sub>2</sub>	-2 + $\gamma \cdot 7.7 = 3.39$

Q	a <sub>1</sub>	a <sub>2</sub>
S <sub>1</sub>	1	3.39
S <sub>2</sub>	7.7	3

Further steps will refine the Q value estimates, and, in the limit, they will converge to their true values.

3. Now let's consider how the Value function changes as the discount factor  $\gamma$  varies between 0 and 1. There are four deterministic policies for this environment, which can be written as  $\pi_{11}$ ,  $\pi_{12}$ ,  $\pi_{21}$  and  $\pi_{22}$ , where  $\pi_{ij}(S_1) = a_i$ ,  $\pi_{ij}(S_2) = a_j$

a. Calculate the value function  $V_{(\gamma)}^{\pi}: S \rightarrow \mathbb{R}$  for each of these four policies (keeping  $\gamma$  as a variable)

$$V_{11}^{\pi}(S_1) = +1 + \gamma V_{11}^{\pi}(S_1), \text{ so } V_{11}^{\pi}(S_1) = 1/(1 - \gamma)$$

$$V_{11}^{\pi}(S_2) = +7 + \gamma V_{11}^{\pi}(S_1) = 7 + \gamma/(1 - \gamma)$$

$$V_{12}^{\pi}(S_1) = V_{11}^{\pi}(S_1) = 1/(1 - \gamma)$$

$$V_{12}^{\pi}(S_2) = 3/(1 - \gamma)$$

$$V_{21}^{\pi}(S_1) = -2 + 7\gamma + \gamma^2 V_{21}^{\pi}(S_1), \text{ so } V_{21}^{\pi}(S_1) = (-2 + 7\gamma)/(1 - \gamma^2)$$

$$V_{21}^{\pi}(S_2) = +7 - 2\gamma + \gamma^2 V_{21}^{\pi}(S_2), \text{ so } V_{21}^{\pi}(S_2) = (7 - 2\gamma)/(1 - \gamma^2)$$

$$V_{22}^{\pi}(S_1) = -2 + 3\gamma/(1 - \gamma)$$

$$V_{22}^{\pi}(S_2) = 3/(1 - \gamma)$$

b. Determine for which range of values of  $\gamma$  each of the policies  $\pi_{11}, \pi_{12}, \pi_{21}, \pi_{22}$  is optimal

$\pi_{11}$  is optimal when

$$0 < V_{11}^{\pi}(S_1) - V_{21}^{\pi}(S_1) = ((1 + \gamma) - (-2 + 7\gamma))/(1 - \gamma^2) = (3 - 6\gamma)/(1 - \gamma^2), \text{ i.e. } 0 \leq \gamma \leq 0.5$$

$\pi_{22}$  is optimal when

$$0 < V_{22}^{\pi}(S_2) - V_{21}^{\pi}(S_2) = (3(1 + \gamma) - (7 - 2\gamma))/(1 - \gamma^2) = (-4 + 5\gamma)/(1 - \gamma^2), \text{ i.e. } 0.8 \leq \gamma < 1.0$$

$\pi_{21}$  is optimal for  $0.5 \leq \gamma \leq 0.8$

$\pi_{12}$  is never optimal because it is dominated by  $\pi_{11}$  when  $\gamma < 2/3$  and by  $\pi_{22}$  when  $\gamma > 0.6$

---