

# PAI exam 2014, Problem 7

Reminder: Q-learning update

$$Q(x, a) \leftarrow (1-\lambda) Q(x, a) + \lambda [r + \gamma \max_{a'} Q(x', a')] \quad (\text{using } \lambda \text{ instead of } \alpha)$$

Episode 1

A1, Right, B1, 0

$$Q(A1, \text{Right}) \leftarrow (1-\lambda) Q(A1, \text{Right}) + \lambda [r + \gamma \max_{a'} Q(B1, a')]$$

$$Q(A1, \text{Right}) \leftarrow 0.5 \times 0 + 0.5 [0 + 0.5 \times 0]$$

$$Q(A1, \text{Right}) \leftarrow 0$$

All  $Q(x, a)$  are initialized to 0, and most rewards are 0. Therefore, most transitions will actually not update the Q values. We will only write down the ones that do a non-zero update.

$$Q(C1, \text{Right}) \leftarrow 0.5 \times 0 + 0.5 [80 + 0.5 \times 0]$$

$$\Rightarrow \boxed{Q(C1, \text{Right}) \leftarrow 40}$$

Episode 2

$$Q(D3, \text{Up}) \leftarrow 0.5 \times 0 + 0.5 [100 + 0.5 \times 0]$$

$$\Rightarrow \boxed{Q(D3, \text{Up}) \leftarrow 50}$$

### Episode 3

$$Q(B3, \text{Left}) \leftarrow 0.5 \times 0 + 0.5 [25 + 0.5 \times 0]$$

$$\Rightarrow Q(B3, \text{Left}) \leftarrow \frac{25}{2}$$

### Episode 4

$$Q(B2, \text{Down}) \leftarrow 0.5 \times 0 + 0.5 [0 + 0.5 \times \overbrace{\max_{a'} Q(B3, a')}^{25/2}]$$

$$\Rightarrow Q(B2, \text{Down}) \leftarrow \frac{25}{8}$$

$$Q(C3, \text{Up}) \leftarrow 0.5 \times 0 + 0.5 [-80 + 0.5 \times 0]$$

$$\Rightarrow Q(C3, \text{Up}) \leftarrow -40$$

Final answers:

$$\begin{aligned} Q(B1, \text{Right}) &= 0 \\ Q(B2, \text{Down}) &= \frac{25}{8} \\ Q(C1, \text{Right}) &= 40 \\ Q(C3, \text{Right}) &= 0 \end{aligned}$$