



### Problem 1.

Consider the following system with two states  $\{A, B\}$  and two possible actions  $\{a^1, a^2\}$ .

The transition probabilities can be expressed as:

$$P(s'|s, a^1) = \begin{cases} 1 & s=A, s'=A \\ 0 & s=A, s'=B \\ 0 & s=B, s'=A \\ 1 & s=B, s'=B \end{cases} \quad P(s'|s, a^2) = \begin{cases} 0 & s=A, s'=A \\ 1 & s=A, s'=B \\ 1 & s=B, s'=A \\ 0 & s=B, s'=B \end{cases}$$

The reward function is as follows:

$\begin{cases} 2 & \text{moving to state } s'=B \\ 0 & \text{moving to state } s'=A \\ -1 & \text{taking action } a^2 \\ 0 & \text{taking action } a^1 \end{cases}$	moving to state $s'=B$ moving to state $s'=A$ taking action $a^2$ taking action $a^1$
---	--

a) Consider the initial Q-value  $\begin{bmatrix} Q(A, a^1) \\ Q(B, a^1) \end{bmatrix} = \begin{bmatrix} 0 \\ -0.1 \end{bmatrix}$  and  $\begin{bmatrix} Q(A, a^2) \\ Q(B, a^2) \end{bmatrix} = \begin{bmatrix} 0.5 \\ 0 \end{bmatrix}$  with  $\gamma=0.9$  and  $\alpha=0.5$ .

Perform Q-Learning Algorithm with  $\epsilon=0$  (greedy Policy), initial state  $s_0=A$  for four steps:  $(s_0=A, a_0)$ ,  $(s_1, a_1)$ ,  $(s_2, a_2)$ ,  $(s_3, a_3)$ ,  $(s_4, a_4)$ . Show all intermediate Q-Values.

b) Compute the final policy  $\pi_A$  and  $\pi_B$ , after all transitions in part (a).

c) If you would use SARSA instead of Q-Learning, would the intermediate Q-Values and final Policy be different from part (a). Justify your answer.  
(no computations needed in this part).

a)

episode 0  
 (A)  $\overset{\text{greedy}}{\underset{\text{---}}{\overset{\text{---}}{a^1 = \arg \max(a)}}} \rightarrow$  (B)

$$Q(A, a^1) = Q(A, a^1) + \frac{\alpha}{0.5} \left[ R_{t+1} + \gamma \max_{a^2} Q(B, a^2) - Q(A, a^1) \right] = 0.75$$

$$\begin{bmatrix} Q(A, a') \\ Q(B, a') \end{bmatrix} = \begin{bmatrix} 0 \\ -0.1 \end{bmatrix} \quad \begin{bmatrix} Q(A, a^2) \\ Q(B, a^2) \end{bmatrix} = \begin{bmatrix} 0.75 \\ 0 \end{bmatrix}$$

closest Deterministic Policy to  $\epsilon$ -greedy

$$\pi(A) = a^2 \quad \pi(B) = a'$$

episode 1

$$(R) \xrightarrow{\substack{a' = a_1 \sim \pi(B) \\ \epsilon\text{-greedy}}} A$$

$$Q(B, a^2) = \frac{Q(B, a^2)}{0} + \frac{\alpha}{0.5} \left[ R_{t+1} + \gamma \max \frac{Q(A, a)}{0.9} - \frac{Q(B, a^2)}{0} \right] = -0.1625$$

$$\begin{bmatrix} Q(A, a') \\ Q(B, a') \end{bmatrix} = \begin{bmatrix} 0 \\ -0.1 \end{bmatrix} \quad \begin{bmatrix} Q(A, a^2) \\ Q(B, a^2) \end{bmatrix} = \begin{bmatrix} 0.75 \\ -0.1625 \end{bmatrix}$$

closest Deterministic Policy to  $\epsilon$ -greedy

$$\pi(A) = a^2 \quad \pi(B) = a'$$

episode 2

$$(A) \xrightarrow{\substack{a' = a_2 \sim \pi(A) \\ \epsilon\text{-greedy}}} B$$

$$Q(A, a^2) = \frac{Q(A, a^2)}{0.75} + \frac{\alpha}{0.5} \left[ R_{t+1} + \gamma \max \frac{Q(B, a)}{0.9} - \frac{Q(A, a^2)}{0.25} \right] = 0.13$$

$$\begin{bmatrix} Q(A, a') \\ Q(B, a') \end{bmatrix} = \begin{bmatrix} 0 \\ -0.1 \end{bmatrix} \begin{bmatrix} Q(A, a^*) \\ Q(B, a^*) \end{bmatrix} = \begin{bmatrix} 0.83 \\ -0.1625 \end{bmatrix}$$

closest Deterministic policy to  $\pi^{\epsilon\text{-greedy}}$

$$\pi(A) = a^* \quad \pi(B) = a'$$

episode 3

$$B \xrightarrow{a' = a_3 \sim \pi(B)^{\epsilon\text{-greedy}}} B$$

$$Q(B, a') = \frac{Q(B, a')}{-0.1} + \frac{\alpha}{0.5} \left[ R_{t+1} + \gamma \max_{a''} Q(B, a'') - \frac{Q(B, a')}{-0.1} \right] = 0.905$$

$$\begin{bmatrix} Q(A, a') \\ Q(B, a') \end{bmatrix} = \begin{bmatrix} 0 \\ 0.905 \end{bmatrix} \begin{bmatrix} Q(A, a^*) \\ Q(B, a^*) \end{bmatrix} = \begin{bmatrix} 0.83 \\ -0.1625 \end{bmatrix}$$

closest Deterministic policy to  $\pi^{\epsilon\text{-greedy}}$

$$\pi(A) = a^* \quad \pi(B) = a'$$

episode 4

$$B \xrightarrow{a' = a_4 \sim \pi(B)^{\epsilon\text{-greedy}}} B$$

$$Q(B, a') = \frac{Q(B, a')}{0.905} + \gamma \left[ R_{t+1} + \frac{\max Q(B, a)}{0.9} - \frac{Q(B, a')}{0.905} \right] = 1.86$$

$$\begin{bmatrix} Q(A, a') \\ Q(B, a') \end{bmatrix} = \begin{bmatrix} 0 \\ 1.86 \end{bmatrix} \begin{bmatrix} Q(A, a^2) \\ Q(B, a^1) \end{bmatrix} = \begin{bmatrix} 0.83 \\ -0.1625 \end{bmatrix}$$

closest Deterministic Policy to  $\epsilon$ -greedy

$$\pi(A) = a^2 \quad \pi(B) = a^1$$

b)  $\pi(A) = \underset{a \in A}{\text{argmax}} Q(A, a) = a^2$

$\pi(B) = \underset{a \in A}{\text{argmax}} Q(B, a) = a^1$

c) It would be same which is because in this question we use  $\epsilon=0$  (greedy policy) all the time, so the  $\max Q(s, a)$  in Q-learning is equal to  $Q(s, a)$  in the SARSA.

## Problem 2.

Consider the following maze problem with 14 states and four actions  $A = \{U, R, D, L\}$

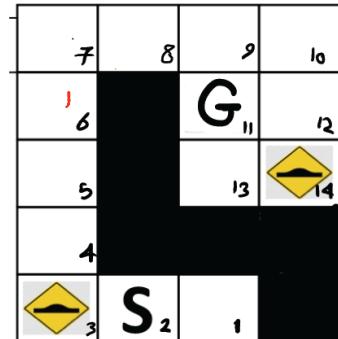
The state transitions are deterministic, and reward for taking any action is -1 and moving to the goal state 100 and bump -10.

a) Set all initial Q-values to zero;  $Q(s, a) = 0$  for all  $s \in S$ ,  $a \in A$ ,  $\alpha = 0.5$  and  $\gamma = 0.9$ .

Run one episode of Q-Learning Algorithm when agent starts from state 5 and follows greedy policy. Note that in the case of equal Q-values, the tie break for actions will be in the following order U, R, D, L. For example

$$\arg\max_{a \in \{U, R, D, L\}} Q(s, a) = \begin{cases} U & \text{if } Q(s, U) = Q(s, R) = Q(s, D) = Q(s, L) \\ R & \text{if } Q(s, R) = Q(s, D) = Q(s, L) > Q(s, U) \\ D & \text{if } Q(s, L) = Q(s, D) > Q(s, U), Q(s, R) \end{cases}$$

b) Show step by step transitions, Q-values and final policy obtained in this single episode.



a)  
 a)  $\underbrace{U=2(5)}_{\epsilon\text{-greedy}} \rightarrow 6$

$$Q(5, U) = \underbrace{Q(5, U)}_0 + \alpha \left[ \underbrace{0.5}_{\text{gamma}} \max \underbrace{\{Q(6, a)\}}_{\text{actions}} - \underbrace{Q(5, U)}_0 \right] = -0.5$$

$$\textcircled{6} \xrightarrow[\epsilon\text{-greedy}]{V = Q(6)} \textcircled{7}$$

$$Q(6, V) = \underbrace{Q(6, U)}_0 + \alpha \left[ \underbrace{R + \delta \max}_{\frac{1}{1-\delta}} (7, A) - \underbrace{Q(6, U)}_0 \right] = -0.5$$

$$\textcircled{7} \xrightarrow[\epsilon\text{-greedy}]{V = Q(7)} \textcircled{7}$$

$$Q(7, U) = \underbrace{Q(7, U)}_0 + \alpha \left[ \underbrace{R + \delta \max}_{\frac{1}{1-\delta}} (7, A) - \underbrace{Q(7, U)}_0 \right] = -0.5$$

$$\textcircled{7} \xrightarrow[\epsilon\text{-greedy}]{R = Q(7)} \textcircled{8}$$

$$Q(7, R) = \underbrace{Q(7, R)}_0 + \alpha \left[ \underbrace{R + \delta \max}_{\frac{1}{1-\delta}} (8, A) - \underbrace{Q(7, R)}_0 \right] = -0.5$$

$$\textcircled{8} \xrightarrow[\epsilon\text{-greedy}]{V = Q(8)} \textcircled{8}$$

$$Q(8, V) = \underbrace{Q(8, U)}_0 + \alpha \left[ \underbrace{R + \delta \max}_{\frac{1}{1-\delta}} (8, A) - \underbrace{Q(8, U)}_0 \right] = -0.5$$

$$\textcircled{8} \xrightarrow[\epsilon\text{-greedy}]{R = Q(8)} \textcircled{9}$$

$$Q(8, R) = \underbrace{Q(8, R)}_0 + \alpha \left[ \underbrace{R + \delta \max}_{\frac{1}{1-\delta}} (9, A) - \underbrace{Q(8, R)}_0 \right] = -0.5$$

$$⑨ \underbrace{U = \mathcal{Q}(q)}_{\epsilon\text{-greedy}} \rightarrow ⑩$$

$$Q(9, U) = \frac{Q(9, U)}{0} + \alpha \left[ R + \delta \max \left( \frac{8}{0.9}, \frac{0}{0} \right) - \frac{Q(9, U)}{0} \right] = -0.5$$

$$⑩ \underbrace{R = \mathcal{Q}(q)}_{\epsilon\text{-greedy}} \rightarrow ⑪$$

$$Q(9, R) = \frac{Q(7, R)}{0} + \alpha \left[ R + \delta \max \left( \frac{8}{0.9}, \frac{0}{0} \right) - \frac{Q(7, R)}{0} \right] = -0.5$$

$$⑪ \underbrace{U = \mathcal{Q}(10)}_{\epsilon\text{-greedy}} \rightarrow ⑫$$

$$Q(10, U) = \frac{Q(10, U)}{0} + \alpha \left[ R + \delta \max \left( \frac{10}{0.9}, \frac{0}{0} \right) - \frac{Q(10, U)}{0} \right] = -0.5$$

$$⑫ \underbrace{R = \mathcal{Q}(10)}_{\epsilon\text{-greedy}} \rightarrow ⑬$$

$$Q(10, R) = \frac{Q(10, R)}{0} + \alpha \left[ R + \delta \max \left( \frac{10}{0.9}, \frac{0}{0} \right) - \frac{Q(10, R)}{0} \right] = -0.5$$

$$⑬ \underbrace{D = \mathcal{Q}(10)}_{\epsilon\text{-greedy}} \rightarrow ⑭$$

$$Q(10, D) = \frac{Q(10, D)}{0} + \alpha \left[ R + \gamma \max_{\text{a}} \left( \frac{Q(12, a)}{0} - \frac{Q(10, D)}{0} \right) \right] = -0.5$$

$$(12) \quad \underbrace{V = Q(12)}_{\epsilon\text{-greedy}} \rightarrow (10)$$

$$Q(12, D) = \frac{Q(12, D)}{0} + \alpha \left[ R + \gamma \max_{\text{a}} \left( \frac{Q(10, a)}{0} - \frac{Q(12, D)}{0} \right) \right] = -0.5$$

$$(10) \quad \underbrace{L = Q(10)}_{\epsilon\text{-greedy}} \rightarrow (9)$$

$$Q(10, L) = \frac{Q(10, L)}{0} + \alpha \left[ R + \gamma \max_{\text{a}} \left( \frac{Q(9, a)}{0} - \frac{Q(10, L)}{0} \right) \right] = -0.5$$

$$(9) \quad \underbrace{P = Q(9)}_{\epsilon\text{-greedy}} \rightarrow (6)$$

$$Q(9, D) = \frac{Q(9, D)}{0} + \alpha \left[ R + \gamma \max_{\text{a}} \left( \frac{Q(6, a)}{0} - \frac{Q(9, D)}{0} \right) \right] = 49.5$$

$$b) \quad \arg\max Q(1, a) = V$$

$$\arg\max Q(2, a) = V$$

$$\arg\max Q(3, a) = V$$

$\operatorname{argmax} Q(4, a) = U$

$\operatorname{argmax} Q(5, a) = R$

$\operatorname{argmax} Q(6, a) = R$

$\operatorname{argmax} Q(7, a) = D$

$\operatorname{argmax} Q(8, a) = D$

$\operatorname{argmax} Q(9, a) = D$

$\operatorname{argmax} Q(10, a) = U$

$\operatorname{argmax} Q(11, a) = U$

$\operatorname{argmax} Q(12, a) = R$

$\operatorname{argmax} Q(13, a) = U$

$\operatorname{argmax} Q(14, a) = U$

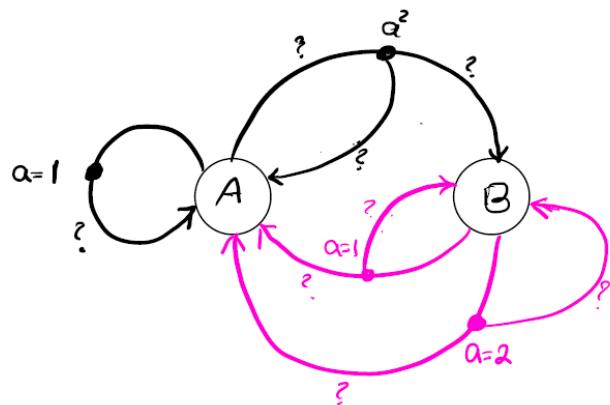
### Problem 3.

Consider the following system with two states  $\{A, B\}$  and two actions  $\{a^1, a^2\}$ . The system state transition is unknown and learning should be achieved through interactions. Consider the following state-action-reward obtained through Softmax Policy in Actor-Critic algorithm.

$$(S_0 = A, a_0 = a^1, r = 10), (S_1 = A, a_1 = a^2, r = -5), (S_2 = B, a_2 = a^1, r = 40),$$

$$(S_3 = A, a_3 = a^2, r = -5), (S_4 = B, a_4 = a^2, r = 20), (S_5 = A, a_5 = a^1, r = +10), S_5 = A$$

Set the initial preferences and state values to zero. Use  $\alpha = 0.5$ ,  $\beta = 0.1$  and  $\gamma = 0.9$  and show all intermediate preferences, state values and policies.



Questions about the HW should be directed to TA, Begum Taskazan, at [taskazan.b@northeastern.edu](mailto:taskazan.b@northeastern.edu).

episode 0

$$\pi(a'|A) = \frac{e^{H(A, a')}}{e^{H(A, a')} + e^{H(A, a')}} = \frac{e^0}{e^0 + e^0} = 0.5 \quad a'$$

$$S_0 = A \xrightarrow{a' = \pi(A)} A \quad \pi(a^t|A) = \frac{e^{H(A, a^t)}}{e^{H(A, a^t)} + e^{H(A, a')}} = \frac{e^0}{e^0 + e^0} = 0.5 \quad a^t$$

$$s_t = r_{t+1} + \gamma V(s_{t+1}) - V(s_t) = 10 + \gamma V(A) - V(A) = 10$$

$$V(A) = V(A) + \alpha s_t = 5$$

$$H(A, a') = H(A, a') + \frac{\beta}{0.1} \times S_t(1 - \pi(a'|A)) = 0.5$$

episode 1

$$\pi(a'|A) = \frac{e^{H(A, a')}}{e^{H(A, a')} + e^{H(A, a')}} = \frac{e^{0.5}}{e^0 + e^{0.5}} = 0.623 \quad a'$$

$$S_1 = A \xrightarrow{a^t = \pi(A)} B \quad \pi(a^t|A) = \frac{e^{H(A, a^t)}}{e^{H(A, a^t)} + e^{H(A, a')}} = \frac{e^0}{e^0 + e^{0.5}} = 0.377 \quad a^t$$

$$s_t = r_{t+1} + \gamma V(s_{t+1}) - V(s_t) = -5 + \gamma V(B) - V(A) = -10$$

$$V(A) = V(A) + \alpha s_t = 0$$

$$H(A, a^t) = H(A, a^t) + \frac{\beta}{0.1} \times S_t(1 - \pi(a^t|A)) = -0.623$$

episode 2

$$\pi(a'|B) = \frac{e^{H(B, a')}}{e^{H(B, a')} + e^{H(B, a')}} = \frac{e^0}{e^0 + e^0} = 0.5 \quad a'$$

$$S_2 = B \xrightarrow{a' = \pi(B)} A \quad \pi(a^t|B) = \frac{e^{H(B, a^t)}}{e^{H(B, a^t)} + e^{H(B, a')}} = \frac{e^0}{e^0 + e^0} = 0.5 \quad a^t$$

$$s_t = r_{t+1} + \gamma V(s_{t+1}) - V(s_t) = 40 + \frac{\gamma}{0.1} V(A) - V(B) = 40$$

$$V(B) = V(B) + \frac{\alpha}{0.5} S_t = 20$$

$$H(B, a') = H(B, a') + \frac{\beta}{0.1} \times S_t (1 - \pi(a' | B)) = 2$$

episode 3

$$S_3 = A \xrightarrow{a^2 = \pi(A)} B$$

$$\begin{aligned} \pi(a' | A) &= \frac{e^{H(A, a')}}{e^{H(A, a')} + e^{H(B, a')}} = \frac{e^{0.75}}{e^{0.75} + e^{-0.623}} = 0.755 \\ \pi(a^2 | A) &= \frac{e^{H(A, a^2)}}{e^{H(A, a^2)} + e^{H(B, a')}} = \frac{e^{-0.623}}{e^{-0.623} + e^{0.75}} = 0.245 \end{aligned}$$

$$S_t = R_{t+1} + \gamma V(S_{t+1}) - V(S_t) = -5 + \frac{\gamma V(B)}{0.9} - \frac{V(A)}{0} = 13$$

$$V(A) = V(A) + \frac{\alpha}{0.5} S_t = 6.5$$

$$H(A, a) = H(A, a^2) + \frac{\beta}{0.1} \times S_t (1 - \pi(a^2 | A)) = 0.3585$$

episode 4

$$S_4 = B \xrightarrow{a^2 = \pi(B)} A$$

$$\begin{aligned} \pi(a' | B) &= \frac{e^{H(B, a')}}{e^{H(B, a')} + e^{H(A, a')}} = \frac{e^2}{e^2 + e^0} = 0.88 | a' \\ \pi(a^2 | B) &= \frac{e^{H(B, a^2)}}{e^{H(B, a^2)} + e^{H(A, a')}} = \frac{e^0}{e^2 + e^0} = 0.119 | a^2 \end{aligned}$$

$$S_t = R_{t+1} + \gamma V(S_{t+1}) - V(S_t) = 20 + \frac{\gamma V(A)}{0.9} - \frac{V(B)}{0} = 5.85$$

$$V(B) = V(B) + \frac{\alpha}{0.5} S_t = 22.925$$

$$H(B, a^2) = H(B, a^2) + \frac{\beta}{0.1} \times S_t (1 - \pi(a^2 | B)) = 0.515$$

episode 5

$$S_t = A \xrightarrow{a^t = \pi(A)} A$$

$$\pi(a^1|A) = \frac{e^{H(A, a^1)}}{e^{H(A, a^1)} + e^{H(A, a^2)}} = \frac{e^{0.5}}{e^{0.5} + e^{0.3585}} = 0.5853 \quad a^1$$

$$\pi(a^2|A) = \frac{e^{H(A, a^2)}}{e^{H(A, a^1)} + e^{H(A, a^2)}} = \frac{e^{0.3585}}{e^{0.5} + e^{0.3585}} = 0.4147 \quad a^2$$

$$S_t = R_{t+1} + \gamma V(S_{t+1}) - V(S_t) = 10 + 0.9V(A) - V(A) = 9.35$$

$$V(A) = V(A) + \underbrace{\alpha}_{0.5} S_t = 11.175$$

$$H(A, a^1) = H(A, a^1) + \underbrace{\beta}_{0.1} \times S_t (-\pi(a^1|A)) = 0.934$$

$$H(A, a^1) = 0.934 \quad H(B, a^1) = 2$$

$$H(A, a^2) = 0.3585 \quad H(B, a^2) = 0.515$$

$$\pi(a^1|A) = \frac{e^{H(A, a^1)}}{e^{H(A, a^1)} + e^{H(A, a^2)}} = 0.64$$

$$\pi(a^2|A) = \frac{e^{H(A, a^2)}}{e^{H(A, a^1)} + e^{H(A, a^2)}} = 0.36$$

$$\pi(a^1|B) = \frac{e^{H(B, a^1)}}{e^{H(B, a^1)} + e^{H(B, a^2)}} = 0.82$$

$$\pi(a^2|B) = \frac{e^{H(B, a^2)}}{e^{H(B, a^1)} + e^{H(B, a^2)}} = 0.18$$