# Midterm Review RL

## Q1. RL

Pacman is in an unknown MDP where there are three states [A, B, C] and two actions [Stop, Go]. We are given the following samples generated from taking actions in the unknown MDP. For the following problems, assume $\gamma = 1$ and $\alpha = 0.5$.

**(a)** We run Q-learning on the following samples:

| s | a | s' | r |
|---|------|---|----|
| A | Go   | B | 2  |
| C | Stop | A | 0  |
| B | Stop | A | -2 |
| B | Go   | C | -6 |
| C | Go   | A | 2  |
| A | Go   | A | -2 |

What are the estimates for the following Q-values as obtained by Q-learning? All Q-values are initialized to 0.

**(i)** $Q(C, Stop) =$ _____

**(ii)** $Q(C, Go) =$ _____

**(b)** For this next part, we will switch to a feature based representation. We will use two features:

- $f_1(s, a) = 1$
- $f_2(s, a) = \begin{cases} 1 & a = \text{Go} \\ -1 & a = \text{Stop} \end{cases}$

Starting from initial weights of 0, compute the updated weights after observing the following samples:

| s | a | s' | r |
|---|------|---|---|
| A | Go   | B | 4 |
| B | Stop | A | 0 |

What are the weights after the first update? (using the first sample)

**(i)** $w_1 =$ _____

**(ii)** $w_2 =$ _____

What are the weights after the second update? (using the second sample)

**(iii)** $w_1 =$ _____

**(iv)** $w_2 =$ _____

# Q2. Q-uagmire

Consider an unknown MDP with three states ($A$, $B$ and $C$) and two actions ($\leftarrow$ and $\rightarrow$). Suppose the agent chooses actions according to some policy $\pi$ in the unknown MDP, collecting a dataset consisting of samples $(s, a, s', r)$ representing taking action $a$ in state $s$ resulting in a transition to state $s'$ and a reward of $r$.

| $s$ | $a$ | $s'$ | $r$ |
|---|---|---|---|
| $A$ | $\rightarrow$ | $B$ | 2 |
| $C$ | $\leftarrow$ | $B$ | 2 |
| $B$ | $\rightarrow$ | $C$ | $-2$ |
| $A$ | $\rightarrow$ | $B$ | 4 |

You may assume a discount factor of $\gamma = 1$.

**(a)** Recall the update function of $Q$-learning is:

$$Q(s_t, a_t) \leftarrow (1 - \alpha)Q(s_t, a_t) + \alpha \left( r_t + \gamma \max_{a'} Q(s_{t+1}, a') \right)$$

Assume that all $Q$-values are initialized to 0, and use a learning rate of $\alpha = \frac{1}{2}$.

**(i)** Run $Q$-learning on the above experience table and fill in the following $Q$-values:

$Q(A, \rightarrow) = $ _____   $Q(B, \rightarrow) = $ _____

**(ii)** After running $Q$-learning and producing the above $Q$-values, you construct a policy $\pi_Q$ that maximizes the $Q$-value in a given state:

$$\pi_Q(s) = \arg\max_a Q(s, a).$$

What are the actions chosen by the policy in states $A$ and $B$?

$\pi_Q(A)$ is equal to:

- ○ $\pi_Q(A) = \leftarrow$.
- ○ $\pi_Q(A) = \rightarrow$.
- ○ $\pi_Q(A) = $ Undefined.

$\pi_Q(B)$ is equal to:

- ○ $\pi_Q(B) = \leftarrow$.
- ○ $\pi_Q(B) = \rightarrow$.
- ○ $\pi_Q(B) = $ Undefined.

**(b)** Use the empirical frequency count model-based reinforcement learning method described in lectures to estimate the transition function $\hat{T}(s, a, s')$ and reward function $\hat{R}(s, a, s')$. (Do not use pseudocounts; if a transition is not observed, it has a count of 0.)

Write down the following quantities. You may write N/A for undefined quantities.

$\hat{T}(A, \rightarrow, B) = $ _____   $\hat{R}(A, \rightarrow, B) = $ _____

$\hat{T}(B, \rightarrow, A) = $ _____   $\hat{R}(B, \rightarrow, A) = $ _____

$\hat{T}(B, \leftarrow, A) = $ _____   $\hat{R}(B, \leftarrow, A) = $ _____

**(c)** This question considers properties of reinforcement learning algorithms for *arbitrary* discrete MDPs; you do not need to refer to the MDP considered in the previous parts.

**(i)** Which of the following methods, at convergence, provide enough information to obtain an optimal policy? (Assume adequate exploration.)

☐ Model-based learning of $T(s, a, s')$ and $R(s, a, s')$.

☐ Direct Evaluation to estimate $V(s)$.

☐ Temporal Difference learning to estimate $V(s)$.

☐ Q-Learning to estimate $Q(s, a)$.

**(ii)** In the limit of infinite timesteps, under which of the following exploration policies is $Q$-learning guaranteed to converge to the optimal Q-values for all state? (You may assume the learning rate $\alpha$ is chosen appropriately, and that the MDP is ergodic: i.e., every state is reachable from every other state with non-zero probability.)

☐ A fixed policy taking actions uniformly at random.

☐ A greedy policy.

☐ An $\epsilon$-greedy policy

☐ A fixed optimal policy.