

## CS 5180 Exercise-2

i) a) State space :-  $\{(x, y)\}$

where,  $0 \leq x \leq 10$ ;  $0 \leq y \leq 10$

Action space :-  $\{\text{UP}, \text{DOWN}, \text{LEFT}, \text{RIGHT}\}$

b) For  $P(s', x | s, a)$

The total no. of possible states are  $x \cdot y$  in which 17 states are not possible. Therefore,

$$\text{no. of possible states} := (11 \times 11) - 17 = 104$$

Total no. of possible actions in each state = 4

The no. of possible next states for a given action = 3

Total no. of possible non-zero rows =  $(104)(4)(3)$

number of different states =  $= 1248$ .

2) a) Episodic task:

We know, total return when the Episode is completed in  $T$  steps, then,

$$G_{t+} = R_t + \gamma R_{t+1} + \gamma^2 R_{t+2} + \dots + \gamma^{T-1} R_T \quad (a)$$

$$= 0 + \gamma(0) + \gamma^2(0) + \dots + \gamma^{T-1}(-1)$$

$$G_{t+} = -\gamma^{T-1}$$

Continuing task:-

Total return  $G_{t+} = R_t + \gamma R_{t+1} + \dots + \gamma^K(-1) + \dots + \gamma^{K+n} R_{t+k+n}$

$$\Rightarrow G_{t+} = 0 + \gamma(0) + \dots + \gamma^K(-1) + \dots + \gamma^{K+n}(-1)$$

$$H = \text{state } G_{t+} = -(\gamma^K + \gamma^{K+1} + \dots)$$

This cannot be simplified further

For episodic tasks there is fixed return value for each episode and for the continuing task there will be no finite value for the return.

3) a) we know,

$$G_t = R_{t+1} + \gamma G_{t+1}$$

$$G_5 = 0$$

$$G_4 = R_5 + \gamma G_5 \Rightarrow 2.$$

$$G_3 = R_4 + \gamma G_4 \Rightarrow 3 + (0.5)(2) = 4$$

$$G_2 = R_3 + \gamma G_3 \Rightarrow 6 + (0.5)(4) = 8$$

$$G_1 = R_2 + \gamma G_2 \Rightarrow 2 + (0.5)(8) = 6$$

$$G_0 = R_1 + \gamma G_1 \Rightarrow (-1) + (0.5)(6) = 2$$

b) Given,  $r = 0.9$   $R_1 = 2$ ;  $R_2 = R_3 = R_4 = \dots = 7$

$$G_{t_0} = 2 + r G_{t_1}$$

$$= 2 + 0.9 [R_2] + (0.9)^2 (R_3) + \dots - (0.9)^\infty (R_\infty)$$

$$= 2 + 0.9(7) + (0.9)^2 (0.7) + \dots - (0.9)^\infty (7)$$

$$= 2 + 7 \left[ 0.9 + (0.9)^2 + (0.9)^3 + \dots \right] \quad \text{(Geometric Progression)}$$

$$= 65$$

$$G_p = R_2 + \gamma R_3 + \gamma^2 R_4 + \dots$$

$$= 7 + \gamma(7) + \gamma^2(7) + \dots$$

$$S \in \mu(2.0) + \delta(-1)R + \delta(-1)$$

$$= 7 \left[ 1 + \gamma + \gamma^2 + \gamma^3 + \dots \right]$$

$$= 7 \left[ \frac{1}{1-\gamma} \right] = 7 \left( \frac{1}{1-0.9} \right)$$

$$S \in 7(2.0) + (-1)(-1)R + \delta(-1) = 14$$

$$= 70$$

u) For 'up' action, totalization will be -1

$$V(\text{UP}) = 50 + \gamma(-1) + \gamma^2(-1) + \dots + \gamma^{100}(-1)$$

$$= 50 - \left[ \gamma + \gamma^2 + \gamma^3 + \dots + \gamma^{100} \right]$$

$$V(\text{UP}) = 50 - \left[ \frac{\gamma(1-\gamma^{100})}{1-\gamma} \right]$$

For 'DOWN' action,

$$v(\text{DOWN}) = -50 + \gamma(1) + \gamma^2(1) + \gamma^3(1) + \dots + \gamma^{100}(1)$$
$$= -50 + \left[ \frac{\gamma(1-\gamma^{100})}{1-\gamma} \right]$$

The agent will choose the action with highest value function.

For choosing 'UP' :-  $v(\text{UP}) > v(\text{DOWN})$

For choosing 'DOWN' :-  $v(\text{UP}) < v(\text{DOWN})$

5) a) The signs of rewards are not important, only the interval between them is important.

$$\text{Eq: } 3.8, G_t = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} R_n$$

Let  $G'_t$  be the new return after adding 'c' to all the reward signs.

$$G'_t = (R_{t+1} + c) + \gamma(R_{t+2} + c) + \gamma^2(R_{t+3} + c) + \dots$$

$$= R_{t+1} + c + \gamma R_{t+2} + \gamma c + \gamma^2 R_{t+3} + \gamma^2 c + \dots$$

$$= c[1 + \gamma + \gamma^2 + \gamma^3 + \dots] + [R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots]$$

$$G_t' = c \left[ \frac{1}{1-r} \right] + G_t$$

$$G_t' = v_c + G_t$$

$$\Rightarrow v_c = c \left[ \frac{1}{1-r} \right]$$

b) for tasks such as maze running, adding a constant for all the rewards may impact the outcomes of the optimal policy.

Let's say hitting the walls result in '-1'

As the agent tries to increase the return this will start avoiding hitting the walls as it will reduce the final return.

Let's say,  $c=2$  is added to all the rewards.

then hitting the wall will result in '+1' reward. The policy will try to increase the final return. therefore it may stuck in hitting the walls for short term rewards instead of finding the exit for the maze.

6) a) We know,

$$V_{\pi}(s) = E_{\pi} \left[ R_{t+1} + \gamma G_{t+1} \mid S_t = s \right]$$

$$\Rightarrow V_{\pi}(l) = 0 + 0.9 \left[ (0.25)(2.3) + (0.25)(0.4) + (0.25)(-0.4) + (0.25)(0.7) \right]$$

$$= 0 + (0.25)(2.3) + (0.25)(0.4) + (0.25)(-0.4) + (0.25)(0.7)$$

$$V_{\pi}(l) = 0.9 \left[ 0.75 \right] = 0.675 \approx 0.7$$

b) We know,

$$V_*(s) = \max_a E_{\pi_*} \left[ R_{t+1} + \gamma G_{t+1} \mid S_t = s, A_t = a \right]$$

$$V_*(l) = 0 + 0.9 \left[ \max(19.8, 19.8, 16, 16) \right]$$

$$= 0.9(19.8) = 17.82 \approx 17.8$$

7) a) As the probability of choosing an action is equal we can guess the  $V(A)$  as 0.5.

$$\begin{aligned} \text{we know, } v(A) &= \pi(L|A) [0 + v(L)] + \pi(R|A) [0 + v(R)] \\ &= (0.5)(0+0) + (0.5)(0+1) \end{aligned}$$

$$V(A) = 0.5(0) + 0.5 \\ = 0.5$$

b) The value function drops as they move away from the right terminal. So,  $V(E) = 0.5$ ,  $V(D) = (0.5)^2$ ,

$$V(C) = 0.5^3, V(B) = 0.5^4, V(A) = 0.5^6$$

$$V(E) = 0.5(V(D) + V(R)) = 0.5(0+1) = 0.5$$

$$V(D) = 0.5(0+0.5^2) = 0.5^2$$

$$V(C) = 0.5(0+0.5^3) = 0.5^3$$

$$V(B) = 0.5(0+0.5^4) = 0.5^4$$

$$V(A) = 0.5(0+0.5^6) = 0.5^5$$

c)  $V(S_n) = 0.5^n$ , where  $S_n$  is the  $n^{th}$  state from R.

8) a) we know, the A is in order and present

How to calculate expectation value in order

$$v(\text{high}) = \pi(\text{search} | \text{high}) \left[ \alpha(x_s + \gamma \cdot v(\text{high})) + (1-\alpha)(x_w + \gamma \cdot v(\text{low})) \right]$$

$$v(\text{low}) = \pi(\text{search} | \text{low}) \left[ (1-\beta)(x_3 + \gamma \cdot v(\text{high})) + \beta(x_s + \gamma \cdot v(\text{low})) \right]$$

$$+ \pi(\text{wait} | \text{low}) \left[ x_w + \gamma \cdot v(\text{low}) \right] + \pi(\text{recharge} | \text{low}) \left[ 0 + \gamma \cdot v(\text{high}) \right]$$

b) by substituting, we get.

$$(0.685) v(\text{high}) = 0.135 \cdot v(\text{low}) + 5$$

$$(0.55) v(\text{low}) = 0.45 \cdot v(\text{high}) + 1.5$$

by solving we get,  $v(\text{high}) = 9.33, v(\text{low}) = 10.36$

c)  $v(\text{high}) = (0.135 \cdot v(\text{low}) + 5) / 0.685$

$$v(\text{low}) = \theta \cdot (3 + 0.9 \cdot v(\text{low})) + (1-\theta) (0.9 \cdot v(\text{high}))$$

$$v(\text{low}) = (6.55 - 3.55(\theta)) / (0.83 + 0.73 \cdot \theta)$$

Choosing the value of  $\theta$  which maximizes the value of  $v(\text{low})$  maximizes  $v(\text{high})$  as well.

$$\text{At } \theta=0, v(\text{low}) = 7.89, v(\text{high}) = 8.85.$$

$$q) \quad a) \quad v_{\pi}(s) = \sum_a \pi(a|s) \cdot q_{\pi}(s, a)$$

$$b) \quad q_{\pi}(s, a) = \sum_{s', r} p(s', r | s, a) [r + \gamma v_{\pi}(s')]$$

$$c) \quad q_{\pi}(s, a) = \sum_{s', r} p(s', r | s, a) \left[ r + \gamma \sum_{a'} \pi(a'|s) q_{\pi}(s', a') \right]$$

$r + \gamma v_{\pi}(s')$  is called the  $\text{Q-value}$ .