Q1. (a).  State space : $S = \{ S_{blank} \}$

Chenghao Wang

Action space : $A = \{ a_{up}, a_{down}, a_{left}, a_{right} \}$

(b) Dynamic of $\{(0,0), DOWN\}$

$$P\{(0,0), 0 \mid (0,0), DOWN\} = 0.9$$
$$P\{(0,1), 0 \mid (0,0), DOWN\} = 0.1$$

Dynamic of $\{(1,5), UP\}$

$$P\{(1,6), 0 \mid (1,5), UP\} = 0.8$$

$$P\{(1,5), 0 \mid (1,5), UP\} = 0.2$$

Dynamic of $\{(9,10), RIGHT\}$

$$P\{(10,10), 1 \mid (9,10), RIGHT\} = 0.8$$
$$P\{(9,10), 0 \mid (9,10), RIGHT\} = 0.1$$
$$P\{(9,9), 0 \mid (9,10), RIGHT\} = 0.1$$

Q2 (a)  $G = R_0 + \gamma R_1 + \gamma^2 R_2 + \cdots + \gamma^T R_T = -\gamma^T$

where T is terminal time step

------------------------------------------------

In continue task, having a discountinue factor could be helpful because it prevents blowing up. But in epsodic task there is no need to having discountinue factor.

(b) Based on the reward setting, the goal for agent here is TO ESCAPING THE MAZE, rather then ESCAPING THE MAZE AS FAST AS POSSIBLE. So there shows no improvement for the training process. If a small negative reward is added to each non terminal step, the effectively communtation could be achieved.

$Q_3$ (a) The sign of those rewards is important when the task is episodic, and the sign become less important for a continuing task.

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \cdots$$

$$\Rightarrow G_{tc} = (R_{t+1} + c) + \gamma(R_{t+2} + c) + \gamma^2(R_{t+3} + c) + \cdots$$

$$= \sum_{k=0}^{\infty} \gamma^k (R_{t+k+1} + c)$$

$$= \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} + \sum_{k=0}^{\infty} \gamma^k c = G_t + \sum_{k=0}^{\infty} \gamma^k c$$

$$V_a(s) = E_a[G_t | S_t = s] = E_a\left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \Big| S_t = s\right]$$

$$V_{ac}(s) = E_a\left[G_t + \sum_{k=0}^{\infty} \gamma^k c \Big| S_t = s\right] = V_a(s) + \sum_{k=0}^{\infty} \gamma^k c$$

$$= V_a(s) + \frac{c}{1-\gamma}$$

(b) This would have effect. For example. a maze. the reward for all non-exit grids is -0.1, the reward for exit is 5. if we add a constant 6 to all the rewards. the agent will never escape the maze.

4.  (a) $V_a(\text{Center}) = P\Big(\big(\text{Center UP}\big), 2.3 \mid \text{Center. UP}\Big) \times \big[0 + 0.9 \times 2.3\big]$

$+ P\Big(\big(\text{Center LEFT}\big), 0.7 \mid \text{Center. LEFT}\Big) \times \big[0 + 0.9 \times 0.7\big]$

$+ P\Big(\big(\text{Center RIGHT}\big), 0.4 \mid \text{Center. RIGHT}\Big) \times \big[0 + 0.9 \times 0.4\big]$

$+ P\Big(\big(\text{Center Down}\big), -0.4 \mid \text{Center. DOWN}\Big) \times \big[0 + 0.9 \times -0.4\big]$

$= 0.25 \times 2.07 + 0.25 \times 0.63 + 0.25 \times 0.36 - 0.25 \times 0.36$

$= 0.6125 + 0.1575$

$\approx 0.7$

(b) $V_*(\text{Center}) = P\Big(\big(\text{Center UP}\big), 19.8 \mid \text{Center. UP}\Big) \times \big[0 + 0.9 \times 19.8\big]$

$+ P\Big(\text{Center LEFT}, 19.8 \mid \text{Center. LEFT}\Big) \times \big[0 + 0.9 \times 19.8$

$= 0.5 \times 0.9 \times 19.8 + 0.5 \times 0.9 \times 19.8$

$= 17.8$

5   (a)

$$V_{\bar{a}}(high) = \bar{a}(search|high) \cdot P(high, r_{search}|high, search) \times [r_{search} + \gamma V_{\bar{a}}(high)]$$

$$+ \bar{a}(search|high) \cdot P(low, r_{search}|high, search) \times [r_{search} + \gamma V_{\bar{a}}(low)]$$

$$+ \bar{a}(wait|high) \cdot P(high, r_{wait}|high, wait) \times [r_{wait} + \gamma V_{\bar{a}}(high)]$$

$$= \bar{a}(search|high) \cdot \alpha \times [r_{search} + \gamma V_{\bar{a}}(high)]$$

$$\bar{a}(search|high) \cdot (1-\alpha) \cdot [r_{search} + \gamma V_{\bar{a}}(low)]$$

$$\bar{a}(wait|high) \cdot [r_{wait} + \gamma V_{\bar{a}}(high)]$$

$$V_{\bar{a}}(low) = \bar{a}(wait|low) \cdot P(low, r_{wait}|low, wait) \times [r_{wait} + \gamma V_{\bar{a}}(low)]$$

$$+ \bar{a}(search|low) \cdot P(low, r_{search}|low, search) \times [r_{search} + \gamma V_{\bar{a}}(low)]$$

$$+ \bar{a}(recharge|low) \cdot P(low, 0|low, recharge) \times [0 + \gamma V_{\bar{a}}(high)]$$

$$+ \bar{a}(search|low) \cdot P(high, r_{search}|low, search) \times [-3 + \gamma V_{\bar{a}}(low)]$$

$$= \bar{a}(wait|low) \cdot [r_{wait} + \gamma V_{\bar{a}}(low)]$$

$$+ \bar{a}(search|low) \cdot \beta [r_{search} + \gamma V_{\bar{a}}(low)]$$

$$+ \bar{a}(recharge|low) \cdot \gamma V_{\bar{a}}(high)$$

$$+ \bar{a}(search|low) \cdot (1-\beta) \cdot [-3 + \gamma V_{\bar{a}}(low)]$$

(b)

$$U_{\tilde{a}}(high) = \tilde{a}(search|high) \cdot \alpha \times [r_{search} + \gamma U_{\tilde{a}}(high)]$$
$$+ \tilde{a}(search|high) \cdot (1-\alpha) \cdot [r_{search} + \gamma U_{\tilde{a}}(low)]$$
$$+ \tilde{a}(wait|high) \cdot [r_{wait} + \gamma U_{\tilde{a}}(high)]$$

$$= 1 \times 0.8 \times [10 + 0.9\, U_{\tilde{a}}(high)]$$
$$+ 1 \times 0.2 \times [10 + 0.9\, U_{\tilde{a}}(low)]$$

$$U_{\tilde{a}}(high) = 8 + 0.72\, U_{\tilde{a}}(high) + 2 + 0.18\, U_{\tilde{a}}(low)$$

$$0.28\, U_{\tilde{a}}(high) = 10 + 0.18\, U_{\tilde{a}}(low) \qquad \text{①}$$

$$U_{\tilde{a}}(low) = \tilde{a}(wait|low) \cdot [r_{wait} + \gamma U_{\tilde{a}}(low)]$$
$$+ \tilde{a}(search|low) \cdot \beta\, [r_{search} + \gamma U_{\tilde{a}}(low)]$$
$$+ \tilde{a}(recharge|low) \cdot \gamma U_{\tilde{a}}(high)$$
$$+ \tilde{a}(search|low) \cdot (1-\beta) \cdot [-3 + \gamma U_{\tilde{a}}(high)]$$

$$= 0.5 \times [3 + 0.9 \times U_{\tilde{a}}(low)]$$
$$+ 0.5 \times 0.9 \times U_{\tilde{a}}(high)$$

$$U_{\tilde{a}}(low) = 1.5 + 0.45\, U_{\tilde{a}}(low) + 0.45\, U_{\tilde{a}}(high)$$

$$0.55\, U_{\tilde{a}}(low) = 1.5 + 0.45\, U_{\tilde{a}}(high)$$

$$11\, U_{\tilde{a}}(low) = 30 + 9\, U_{\tilde{a}}(high) \quad \text{②} \qquad \Rightarrow 77 U_{low} = 210 + 63 U_{high}$$

Combine ① and ② 

$$\Rightarrow 14\, U_{\tilde{a}}(high) = 500 + 9\, U_{\tilde{a}}(low) \quad \Rightarrow 63 U_{high} = 2250 + 40.5 U_{low}$$

$$36.5\, U_{low} = 2460$$

$\Rightarrow U_{low} = 67$

$\Rightarrow U_{high} = 79$

$14 U_{high} = 500 + 603$

6 (a) $\quad U_{\bar{a}} = E_{\bar{a}}[G_t | S_t = s]$

$\qquad = \sum_a E_{\bar{a}}[G_t | S_t = s, A_t = a] P[A_t = a | S_t = s]$

$\qquad = \sum_a q_{\bar{a}}(s, a) \bar{\pi}(a|s)$

(b) $q(s,a) = E_{\bar{a}}[G_t | S_t = s, A_t = a]$

$\qquad = \sum_{s'} [G_t | S_t = s, A_t = a, S_{t+1} = s'] P[S_{t+1} = s' | S_t = s, A_t = a]$

$\qquad = \sum_{s'} [G_t | S_{t+1} = s'] P[S_{t+1} = s' | S_t = s, A_t = a]$

$\qquad = \sum_{s,r} \{ r + \gamma E_{\bar{a}}[G_{t+1} | S_{t+1} = s'] \} P(s', r | s a)$

$\qquad = \sum_{s,r} [r + \gamma U_{\bar{a}}(s')] P(s', r | s a)$

(c) $q_{\bar{a}}(s, a) = \sum_{s,r} P(s', r | s, a) [r + \gamma U_{\bar{a}}(s')]$

$\qquad = \sum_{s,r} P(s, r | s a) [r + \sum_{a \in A} \pi(a'|s') \cdot \gamma q_{\bar{a}}(s', a')]$