

Q1)

$$a) S = \{ (x, y) \in \mathbb{N} \times \mathbb{N} \mid 0 \leq x \leq 10, 0 \leq y \leq 1 \}$$

We also consider walls states, which could not be accessible

$$A = \{ \text{LEFT, DOWN, RIGHT, UP} \}$$

b) We have 11x11 states which for each state we have 4 actions (possible action) and for each action we usually

have 3 possible transfer. of course we could not to transfer to wall states and never have transformation from wall states so approximately we will have

$$[11 \times 11 - \text{num. of wall states}] \times 4 \times 3 - \text{num. of wall states}$$

$$= [11 \times 11 - 17] \times 4 \times 3 - 17 = 1231 \text{ non-zero Rows}$$

Q2)-

a-) episodic task with discount

$$G_t = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{T-t-1} R_T + \sum_{K=0}^{T-t-1} \gamma^K R_{t+K+1}$$

$$\text{If } R_{t+1} = \dots = R_{T-1} = 0, R_T = -1$$

$$\Rightarrow G_t = -\gamma^{T-t-1}$$

while for continuing task, the return is $-\gamma^K$

where K is the number of time steps before failure.

So for episodic $G_t = -\gamma^{T-t-1} = -\gamma^{K-1}$, considering

t is K steps before T , failure.

b) Because there is no discount, we do not force the robot to exit the maze as soon as possible. In this case when the robot exits get rewards 1 otherwise it gets zero which does not force it to exit quickly. So we have not communicated effectively to the robot. For improving it'd better to have discount or to consider negative rewards for non-goal states.

So we have not communicated effectively to the robot. For improving it'd better to have discount or to consider negative rewards for non-goal states

3)

a) In continuing task, the sign of rewards is not important and the intervals between them is important, as shown in following. If we add a positive value to all rewards finally the relative values between states will not change

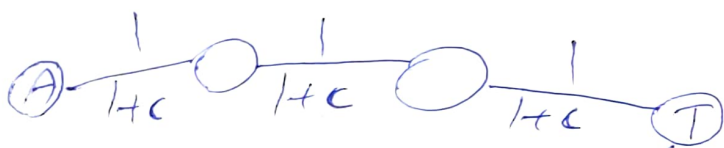
$$G_t(S) = R_{t+1} + \gamma R_{t+2} + \dots = \sum_{K=0}^{\infty} \gamma^K R_{t+K+1} \quad \forall S$$

$$G'_t(S) = (R_{t+1} + c) + \gamma [R_{t+2} + c] + \dots = c \sum_{K=0}^{\infty} \gamma^K + \sum_{K=0}^{\infty} \gamma^K R_{t+K+1} \quad \forall S$$

$$G'_t(S) = c \sum_K \gamma^K + G_t(S) \quad \forall S \Rightarrow V_c = c \sum_K \gamma^K \quad \forall S$$

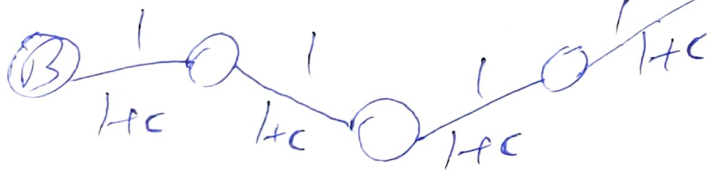
$$\text{If } \gamma < 1 \Rightarrow V_c = \frac{c}{1-\gamma} \quad \forall S$$

b) In the episodic task, adding a constant would be problematic because at first it will change the relative values between states and second it cause some negative rewards to be positive that change our problem.



$$G_+(A) = 3$$

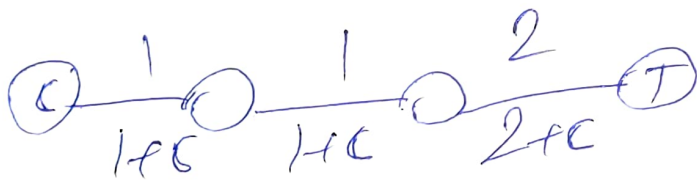
$$G'_+(A) = 3 + 3c$$



$$G_+(B) = 4$$

$$G'_+(B) = 4 + 4c$$

T: is terminal state



$$G_+(C) = 4$$

$$G'_+(C) = 4 + 3c$$

$$G_+(A) - G_+(B) = -1$$

$$G'_+(A) - G'_+(B) = -(1+c)$$

$$G_+(A) - G_+(C) = -1$$

$$G'_+(A) - G'_+(C) = -(1)$$

before adding c to each reward the relative

value between A-B and A-C was same

but after adding c this relative value changes.

which makes our problem distinct from the first one

$$4) V_H(s) \sum \pi(a|s) \sum_{r \in \mathcal{R}} P(s', r | s, a) [1 + \gamma V_H(s')] \\ \forall s \in \mathcal{S}$$

$$\gamma = 0.9$$

$$P(s', r | s, a) = 1 \quad \forall s, s', a, r$$

$$\gamma = 0$$

$$\Rightarrow a) 0.7 = \frac{1}{4} [0 + 0.9 \times 2.3] + \frac{1}{4} [0 + 0.9 \times 0.4] + \frac{1}{4} [0 + 0.9 \times 0.7] \\ + \frac{1}{4} [0 - 0.9 \times 0.4] = 0.675 \sim 0.7$$

$$b) 17.8 = \frac{1}{2} [0 + 0.9 \times 19.8] + \frac{1}{2} [0 + 0.9 \times 19.8] = 17.82$$

5)

a) Guess: $V_{\pi}(A) = \frac{1}{2}$ $V(L) = 0$ $V(R) = 1$

$$V_{\pi}(S) = \sum \pi(a|S) \sum P(S', r | S, a) [r + \gamma V_{\pi}(S')]$$

$\gamma = 1$

$$P(S', r | S, a) = 1$$

$$V_{\pi}(S=R) = 1$$

$$V_{\pi}(S=L) = 0$$

$$r = 0$$

$$V_{\pi}(S=A) = \pi(\text{Right} | S=A) [0 + V_{\pi}(S=R)] + \pi(\text{Left} | S=A) [0 + V_{\pi}(S=L)]$$

$$= \frac{1}{2} \times 1 + \frac{1}{2} \times 0 = \frac{1}{2}$$

$$V_{\pi}(S=R) = 1 \quad V_{\pi}(S=L) = 0$$

b) because there is not discount \Rightarrow Guess:

$$V(L) = 0 \quad V(A) = \frac{1}{6} \quad V(B) = \frac{2}{6} \quad V(C) = \frac{3}{6}$$

$$V(D) = \frac{4}{6} \quad V(E) = \frac{5}{6} \quad V(R) = 1$$

$$\begin{aligned} \text{ex: } V_{\pi}(S \leq E) &= \pi(\text{Right} | S \leq E) [0 + V_{\pi}(S \leq R)] + \\ &\quad \pi(\text{Left} | S \leq E) [0 + V_{\pi}(S \leq D)] = \\ &\quad \frac{1}{2} \times 1 + \frac{1}{2} \times \frac{4}{6} = \frac{5}{6} \end{aligned}$$

c) based on part c we could guess that
for n states the Value Function of states
would be



$$V_{\pi}(S \leq S_1) = \frac{1}{n} \quad V_{\pi}(S \leq S_2) = \frac{2}{n} \quad \dots \quad V_{\pi}(S \leq S_{n-2}) = \frac{n-1}{n}$$

$$V_{\pi}(S \leq L) = 0 \quad V_{\pi}(S \leq R) = 1$$

6) H: high L: low

$$a) V_{\pi}(s) = \sum \pi(a|s) \sum P(s'|a, s_{\text{net}}) [r + \gamma V_{\pi}(s')]$$

$$V_{\pi}(s = \text{High}) = \pi(\text{Search}|H) (\alpha [r_{\text{Search}} + \gamma V_{\pi}(H)] + (1-\alpha) [r_{\text{Search}} + \gamma V_{\pi}(\text{Low})]) \\ + \pi(\text{wait}|H) (\gamma [r_{\text{wait}} + \gamma V_{\pi}(H)])$$

$$V_{\pi}(s = L) = \pi(\text{Search}|L) (\beta [r_{\text{Search}} + \gamma V_{\pi}(L)] + (1-\beta) [-3 + \gamma V_{\pi}(H)]) \\ + \pi(\text{wait}|L) (\gamma [r_{\text{wait}} + \gamma V_{\pi}(L)])$$

$$+ \pi(\text{Recharge}|L) (\gamma [0 + \gamma V_{\pi}(H)])$$

$$b) \alpha = 0.8 \quad \beta = 0.6 \quad \gamma = 0.9 \quad V_{\text{Search}} = 10$$

$$V_{\text{wait}} = 3 \quad \pi(\text{Search} | H) = 1 \quad \pi(\text{wait} | L) = 0.5$$

$$\pi(\text{recharge} | L) = 0.5$$

$$V_{\pi}(H) = 10 + 0.72 V_{\pi}(H) + 0.18 V_{\pi}(L)$$

$$V_{\pi}(L) = 1.5 + 0.45 V_{\pi}(L) + 0.45 V_{\pi}(H)$$

$$\Rightarrow V_{\pi}(L) = 87.7 \quad V_{\pi}(H) = 79.0$$

7)

$$a) V_{\pi}(s) = \sum_a \pi(a|s) q_{\pi}(a|s)$$

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

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