

3.6) ~~Reward~~ for a continuous task $(R) = r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots$
Return

" " " episodic task $(R) = r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots$

In this case since we also introduce discounts in the episodic task, the ~~reward~~ return will be the same as a continuous task:

$$R = r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots$$

Return at each time will be 0, except for a failure where it will be $= -\gamma^{k-1}$

3.7) No, we have not introduced any negative reward on taking a wrong move, meaning a short path (5 grids) and a longer path (8 grids) will give it the same reward (+1) on escape, so the agent doesn't learn anything.

By introducing negative reward (say -1), the rewards will be -4 and -7 respectively instead, according to my examples. So it receives negative reinforcement, such that it learns from the policy / path it took to escape, and can learn and differentiate between good and better paths.

3.14) Figure 3.2

Belman equation:-

$$V_{\pi}(s) = \sum_a \pi(a|s) \sum_{s', r} p(s', r | s, a) [r + \gamma V_{\pi}(s')], \text{ for all } s \in S$$

Here, given:-

$$p(s', r | s, a) = 1 \quad (\text{transition probability})$$

$$\gamma = 0.9$$

$$r = 0$$

$$\pi(a|s) = 1/4 = 0.25 \quad (4 \text{ choices})$$

Therefore:-

$$\begin{aligned} V_{\pi}(0.7) &= 0.25 \times 1 (0 + (0.9 \times 0.2)) \\ &\quad + 0.25 \times 1 (0 + (0.9 \times 0.4)) \\ &\quad + 0.25 \times 1 (0 + (0.9 \times -0.4)) \\ &\quad + 0.25 \times 1 (0 + (0.9 \times 0.7)) \end{aligned}$$

$$= 0.25 \left(\overset{2.07}{0.18} + 0.36 - 0.36 + 0.63 \right)$$

$$= 0.675 \approx 0.7$$

3.22) ~~2~~ Bellman equation:-

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

In this case:-

$$\pi(a|s) = 0.5 \quad (\text{two choices})$$

$$p(s', r | s, a) = 1$$

for $\gamma = 0$:

$$~~V_{\pi}(s) = 0.5 \times 1~~$$

$$V^*(s) = \max(Q_L, Q_R) \quad , \text{for state } s, L = \text{left policy} \\ R = \text{right policy}$$

where:-

$$Q_L = \sum p(s', r | s, a) [r + \gamma V_{\pi}(s')] \\ = 1 [1 + 0 \times 0] \\ = \underline{1}$$

$$Q_R = 1 [0 + 0 \times 2] \\ = 0$$

$$\therefore V^*(s) \text{ for } \gamma = 0 = \underline{\underline{Q_L}}$$

Similarly for $\gamma = 0.9$ & $\gamma = 0.5$:

γ	Q_L	Q_R	
0	①	0	→ Optimal policy = left path
0.9	1	①.8	→ Optimal policy = right path
0.5	①	①	→ Optimal policy = either path is fine

3.23) Bellman equations for Q^* :

Given 2 states → low, high
and 3 actions → search, wait, recharge

$$Q^*(\text{high, search}) = \alpha(R_{\text{search}} + \gamma V^*(\text{high})) + (1-\alpha)(R_{\text{search}} + \gamma V^*(\text{low})) \quad - (1)$$

$$Q^*(\text{high, wait}) = R_{\text{wait}} + \gamma V^*(\text{high}) \quad - (2)$$

$$Q^*(\text{low, search}) = \beta(R_{\text{search}} + \gamma V^*(\text{low})) + (1-\beta)(-3 + \gamma V^*(\text{high})) \quad - (3)$$

$$Q^*(\text{low, recharge}) = \gamma V^*(\text{high}) \quad - (4)$$

$$Q^*(\text{low, wait}) = R_{\text{wait}} + \gamma V^*(\text{low}) \quad - (5)$$

3.25) V^* in terms of Q^* :-

$$V^*(high) = \max(Q^*(high, search), Q^*(high, wait)).$$

$$V^*(low) = \max(Q^*(low, search), Q^*(low, wait), Q^*(low, recharge)).$$