

Solⁿ of Chapter 3

3.1 i) Stock Trader :-

State :- Current price, Chart of previous days,
Amount to be invested

Action :- Buy, sell, hold.

Reward :- ~~∑~~ Total profit = Reward.

ii) Rubix Cube Solver :-

State :- Current Cube State.

Action :- All Valid Cube moves.

Reward :- +100 on correct solve

iii) QWOP game agent:

State :- Players Stance

Action :- Q, W, O, P

Reward :- Distance Covered

3.2 Nope, MDP framework doesn't represent all tasks.

Ex :- Chat Bot

Next response depend on entire chat history
& not just most recent

3.3 There is no fundamental reason to prefer one location over other, it's a free choice.

3.4

Q3.5 $\sum_{s' \in S} \sum_{r \in R} p(s', r | s, a) = 1 \quad \forall s \in S, a \in A(s)$
 Modified eqⁿ ... eq(3.3)

$$\sum_{s' \in S} \sum_{r \in R} p(s', r | s, a) = 1 \quad \forall s \in S, a \in A(s)$$

3.6 Return = $- \gamma^{T-1}$

3.7 $\gamma = 1$ for every episode irrespective of Time taken for each episode
 Hence Discounting should be used

3.8

$$\begin{aligned} R_4 &= R_5 + 0.5 R_5 = 2 \\ R_3 &= R_4 + 0.5 R_4 = 4 \\ R_2 &= R_3 + 0.5 R_3 = 8 \\ R_1 &= R_2 + 0.5 R_2 = 6 \\ R_0 &= R_1 + 0.5 R_1 = 2 \end{aligned}$$

3.9

$$\begin{aligned} R_0 &= 2 + (0.9)7 + (0.9)^2 7 + \dots \\ &= 2 + 7 \times 0.9 [1 + 0.9 + 0.9^2 + \dots] \\ &= 2 + 7 \times 0.9 \left[\frac{1}{1-0.9} \right] \end{aligned}$$

$$R_0 = 2 + \frac{7 \times 0.9}{0.1} = 65$$

$$\begin{aligned} R_1 &= 7 + (0.9)7 + \dots \\ &= 7 [1 + 0.9 + (0.9)^2 + \dots] \\ R_1 &= \frac{7 [1]}{0.1} = 70 \end{aligned}$$

3.10 Summation of $R_i P_i$.

$$3.11. E[R_{t+1} | S_t] = \sum_{a \in A} \pi(a | S_t) \cdot r(a, S_t).$$

$$= \sum_{a \in A} \pi(a, S_t) \cdot \sum_{r \in R} \sum_{S' \in S} r \cdot p(S', r | S_t, a)$$

$$3.12. V_{\pi}(S) = \sum_{a \in A} q_{\pi}(a, S) \cdot \pi(a | S)$$

$$3.13. q_{\pi}(S, a) = E[R_{t+1} + \gamma V_{\pi}(S_{t+1}) | S_t = S, A_t = a]$$

$$= \sum_{S', r} p(S', r | S, a) \cdot E_{\pi}[R_{t+1} + \gamma V_{\pi}(S_{t+1}) | S_{t+1} = S', R_{t+1} = r]$$

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

3.14

		2.3	up
0.7	0.7	0.4	right
left	-0.4		down

$$V(\text{up}) =$$

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

$$= 0.9 \frac{2.3 + 0.7 + 0.4 + \cancel{0} - 0.4}{4}$$

$$= 0.9 \times 0.75 = 0.675 \approx 0.7$$

$$3.15. Q_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$$

$$Q_{t_{\text{new}}}^* = Q_t + \sum_{k=0}^{\infty} \gamma^k C = Q_t + \frac{C}{1-\gamma}$$

$$\therefore E[Q_{t_{\text{new}}} | S_t = S] = E\left[Q_t + \frac{C}{1-\gamma} \mid S_t = S\right]$$

$$V_{\pi_{\text{new}}}(S) = V_{\pi}(S) + \frac{C}{1-\gamma}$$

$$3.25 \quad V_*(s) = \max_a q_*(s, a)$$

$$3.26 \quad q_*(a, s) = \sum_{s', r} p(s', r | s, a) [r + \gamma V_*(s')]$$

$$3.27 \quad \pi^*(s) = \arg \max_a q_*(s, a)$$

$$3.28 \quad \pi^*(s) = \arg \max_a \left[\sum_{s', r} p(s', r | s, a) [r + \gamma V_*(s')] \right]$$

3.29