

Q3.11

$$E[R_{t+1} | s_t = s] = \sum_a \pi(a|s) \sum_{s', r} p(s', r | s, a) r$$

Q3.12

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

Q3.13

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

Q3.14

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

$$V_{center} = \frac{1}{4} \times 0.9 \times (2.3 + 0.7 - 0.4 + 0.4) = 0.675 \approx 0.7$$

Q3.15

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

Adding a Constant C

~~$$G_t^* = G_t + C = C$$~~

$$G_t^* = G_t + \sum_{k=0}^{\infty} \gamma^k C = G_t + \frac{C}{1-\gamma}$$

$$v_{\pi}^*(s) = E[G_t^* | s_t = s] = E[G_t + \frac{C}{1-\gamma} | s_t = s] = E[G_t | s_t = s] + \frac{C}{1-\gamma}$$

The new $v_{\pi}^*(s)$ does not affect the relative difference among states.

Q3.16

The sign of the reward has critical influence on the episodic rewards because episodic tasks use negative rewards to accelerate the agent finishing the task. Thus sign of the agent would impact how the agent moves. Furthermore, if the negative rewards ~~remains~~ negative but the value of it shrinks too much, it will give a wrong signal to the agent that the time of completing the job is not important.

Q3.17

$$\begin{aligned} q_{\pi}(s, a) &= E_{\pi} [G_t | S_t = s, A_t = a] \\ &= E_{\pi} [R_{t+1} + \gamma G_{t+1} | S_t = s, A_t = a] \\ &= \sum_{s', r} p(s', r | s, a) [r + \gamma \sum_{a'} \pi(a', s') q_{\pi}(s', a')] \end{aligned}$$

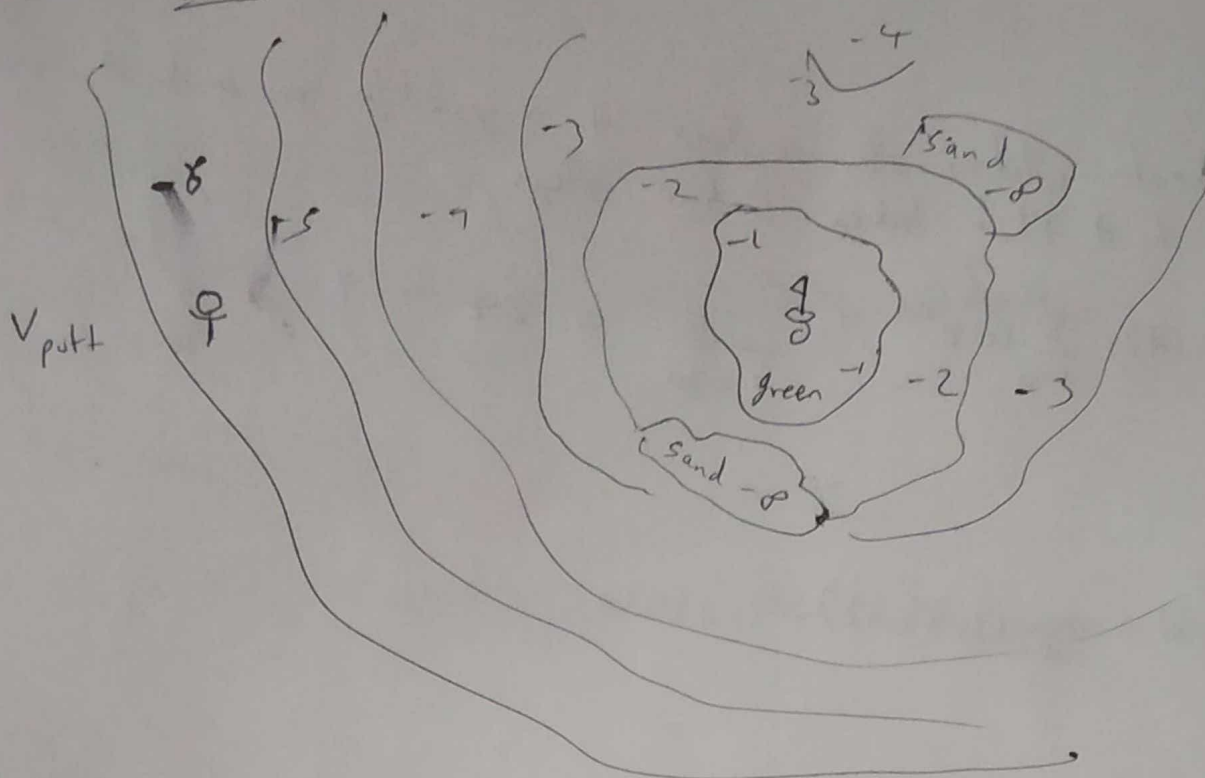
Q3.18

$$\begin{aligned} v_{\pi}(s) &= E_{\pi} [q_{\pi}(s_t, A_t) | S_t = s, A_t = a] \\ &= \sum_a \pi(a | s) q_{\pi}(s, a) \end{aligned}$$

Q3.19

$$\begin{aligned} q_{\pi}(s, a) &= E_{\pi} [R_{t+1} + v_{\pi}(s') | S_t = s, A_t = a] \\ &= \sum_{s', r} p(s', r | s, a) [r + v_{\pi}(s')] \end{aligned}$$

Q.3.20 & 3.21



Optimal state values according to driver when off green, the according to putter on the green.

Optimal policy is to use driver when off green & putter when on green.

Q.3.22

$$G_{n, \text{left}} = \sum_{i=0}^{\infty} \gamma^i = \frac{1}{1-\gamma^2} \quad G_{n, \text{right}} = \sum_{i=0}^{\infty} 2\gamma^{1+2i} = \frac{2\gamma}{1-\gamma^2}$$

If $\gamma > 0.5$, right is optimal

If $\gamma < 0.5$, left is optimal

If $\gamma = 0.5$, both are optimal

Q.3.24

The best solution after reaching A is quickly go back A after moving to A. That takes 5 time steps.

$$v_0(A) = \sum_{t=0}^{\infty} 10\gamma^{5t} = \frac{10}{1-\gamma^5} \approx 24.419$$

Q.3.25

$$v_*(s) = \max_a (q_*(s, a))$$

Q.3.26

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

Q.3.27

$$a_* = \arg \max_a \pi_*(a_*, s) = \arg \max_a q_*(s, a)$$

Q.3.28

$$a_* = \arg \max_a \pi_*(a_*, s) = \arg \max_a \sum_{s', r} p'(s', r | s, a) [r + \gamma v_*(s')]$$

8.3.29

$$\begin{aligned} V_{\pi}(s) &= E_{\pi}[G_t | s_t = s] \\ &= \sum_a [r(s, a) + \gamma \sum_{s'} p(s' | s, a) V_{\pi}(s')] \pi(s, a) \end{aligned}$$

$$\begin{aligned} V_{*}(s) &= E_{\pi}[G_t | s_t = s] \\ &= \sum [r(s, a) + \gamma \sum_{s'} p(s' | s, a) V_{*}(s')] \pi_{*}(s, a) \end{aligned}$$

$$\begin{aligned} q_{\pi}(s, a) &= E_{\pi}[G_t | s_t = s, A_t = a] \\ &= E_{\pi}[R_{t+1} + \gamma G_{t+1} | s_{t+1} = s', A_t = a] \\ &= ~~r(s, a)~~ r(s, a) + \gamma \sum_{s'} p(s' | s, a) \sum_{a'} q_{\pi}(a', s') \pi(a' | s') \end{aligned}$$

$$\begin{aligned} q_{*}(s, a) &= E_{\pi_{*}}[G_t | s_t = s, A_t = a] \\ &= E_{\pi_{*}}[R_{t+1} + \gamma G_{t+1} | s_{t+1} = s', A_t = a] \\ &= r(s, a) + \gamma \sum_{s'} p(s' | s, a) \sum_{a'} q_{*}(a', s') \pi_{*}(a' | s') \end{aligned}$$