

Ques 1. We know  $\pi_{search}$

- For  $s = \text{high}$ ,  $a = \text{search}$ ,  $s' = \text{high}$

$$\pi_{search} = \frac{(1) P(s', \pi=1 | s, a)}{P(s' | s, a)}$$

$$P(s', \pi=1 | s, a) = \alpha \pi_{search}$$

$$P(s' | s, a) = P(s', \pi=1 | s, a) + P(s', \pi=0 | s, a)$$

$$\Rightarrow P(s', \pi=0 | s, a) = 1 - \alpha \pi_{search}$$

- Similarly for  $s = \text{high}$ ,  $a = \text{search}$ ,  $s' = \text{low}$

$$P(s', \pi=1 | s, a) = (1-\alpha) \pi_{search}$$

$$P(s', \pi=0 | s, a) = (1-\alpha) - (1-\alpha) \pi_{search}$$

- Similarly for  $s = \text{low}$ ,  $a = \text{search}$ ,  $s' = \text{low}$

$$P(s', \pi=1 | s, a) = \beta \pi_{search}$$

$$P(s', \pi=0 | s, a) = \beta - \beta \pi_{search}$$

- For  $s = \text{high}$ ,  $a = \text{wait}$ ,  $s' = \text{high}$

$$\pi_{wait} = P(s', \pi=1 | s, a)$$

$$\Rightarrow P(s', \pi=0 | s, a) = 1 - \pi_{wait}$$

- Similarly for  $s = \text{low}$ ,  $a = \text{wait}$ ,  $s' = \text{low}$

$$P(s', \pi=1 | s, a) = \pi_{wait}$$

$$P(s', \pi=0 | s, a) = 1 - \pi_{wait}$$

Table:-

| s    | a      | s'   | $\pi$ | $P(s', \pi   s, a)$                    |
|------|--------|------|-------|--|
| High | Search | High | 0     | $1 - \alpha \pi_{search}$              |
| High | Search | High | 1     | $\alpha \pi_{search}$                  |
| High | Search | Low  | 0     | $(1-\alpha) - (1-\alpha) \pi_{search}$ |
| High | Search | Low  | 1     | $(1-\alpha) \pi_{search}$              |
| Low  | Search | Low  | 0     | $\beta - \beta \pi_{search}$           |
| Low  | Search | Low  | 1     | $\beta \pi_{search}$                   |
| Low  | Search | High | -3    | $1 - \beta$                            |
| Low  | Wait   | Low  | 0     | $1 - \pi_{wait}$                       |
| Low  | Wait   | Low  | 1     | $\pi_{wait}$                           |

| $s$  | $a$      | $s'$ | $n$ | $p(s', n   s, a)$ |
|------|----------|------|-----|-------------------|
| High | wait     | High | 0   | $1 - \pi_{wait}$  |
| High | wait     | High | 1   | $\pi_{wait}$      |
| Low  | Recharge | High | 0   | 1                 |

Ques 3.

(a) Exercise 3.15

The signs of the rewards are not important but intervals between them are.

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

If we add constant to all  $\sum_{k=0}^{\infty}$  rewards then,

$$G_t' = \sum_{k=0}^{\infty} \gamma^k (R_{t+k+1} + c) = G_t + \sum_{k=0}^{\infty} \gamma^k c$$

$$\Rightarrow G_t' = G_t + \frac{c}{1-\gamma} = G_t + v_c$$

$$V_{\pi}(s)' = E_{\pi} [G_t' | S_t = s]$$

$$= E_{\pi} [G_t + v_c | S_t = s]$$

$$= E_{\pi} [G_t | S_t = s] + E_{\pi} [v_c | S_t = s]$$

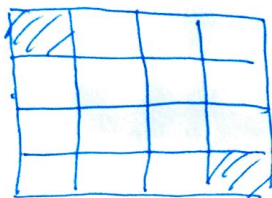
$$= V_{\pi}(s) + v_c$$

$\Rightarrow v_c = \frac{c}{1-\gamma}$   $\left\{ \begin{array}{l} \text{Thus only constant } v_c \text{ is added to } v_{\pi} \text{ on} \\ \text{increasing rewards to } v_{\pi} \end{array} \right\}$

(b) Exercise 3.16

Adding a constant ( $c$ ) to all rewards in an episodic task will effect the tasks.

Consider the following episodic task:-



$(0,0)$  and  $(3,3)$  are terminal states  
reward for each action  $(-1)$

Optimal Policy would ensure shortest path to either of the terminal states.

Now if the rewards are  $-1 + c$  ( $c > 1$ ) then optimal policy changes, also value functions will for any policy will be different as it would never try to go to terminal states for maximizing reward

Ques 5.  $V_{\pi}(s) = \max_{a \in A(s)} q_{\pi}(a, s)$

$$= \max_a E_{\pi} [G_t | A_t = a, S_t = s]$$

$$= \max_a E_{\pi} [R_t + \gamma G_{t+1} | A_t = a, S_t = s]$$

$$= \max_a E_{\pi} [R_t + \gamma V_{\pi}(S_{t+1}) | A_t = a, S_t = s]$$

$$= \max_a \sum_{s'} \sum_{a'} P(s', a' | s, a) \{r + \gamma V_{\pi}(s')\}$$

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