**Q1**

| S | a | s' | r | $P(s', r \mid s, a)$ |
|---|---|---|---|---|
| high | search | high | r | $(\alpha) \, P_{search}(r)$ |
| high | search | low | r | $(1-\alpha) \, P_{search}(r)$ |
| low | search | high | -3 | $(1-\beta) \, P_{search}(r)$ |
| low | search | low | r | $(\beta) \, P_{search}(r)$ |
| high | wait | high | r | $(1) \, P_{wait}(r)$ |
| low | wait | low | r | $(1) \, P_{wait}(r)$ |
| low | recharge | high | 0 | $(1) \, (1)$ |

Here $P_{search}(r)$ is a probability distribution with mean '$r_{search}$'
and $P_{wait}(r)$ is a probability distribution with mean '$r_{wait}$'

$$\Rightarrow \quad E[P_{search}(r)] = \text{'}r_{search}\text{'}$$
$$\text{and} \quad E[P_{wait}(r)] = \text{'}r_{wait}\text{'}$$

Signs of rewards are NOT important and only the intervals between them are. This is because if a large constant is added or subtracted, then all rewards can be made of the same sign. Since the relative order is preserved, the algorithm is not affected. This can be seen below :

we know that

$$V_\pi(s) = E\left[G_t \mid S_t = s\right]$$

where

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

Now, if a constant 'c' is added to all rewards

$\hat{G}_t$ becomes $\sum_{k=0}^{\infty} \gamma^k (R_{t+k+1} + c)$

$$\Rightarrow \hat{G}_t = G_t + \sum_{k=0}^{\infty} c\gamma^k$$

$$\therefore \hat{G}_t = G_t + \frac{c}{1-\gamma} \quad (\text{since } \gamma < 1)$$

Now, $\hat{V_\pi}(s) = E\left[\hat{G}_t \mid S_t = s\right]$

$$= E\left[G_t + \frac{c}{1-\gamma} \mid S_t = s\right]$$

$$= E\left[G_t \mid S_t = s\right] + \frac{c}{1-\gamma}$$

or $\boxed{\hat{V_\pi}(s) = V_\pi(s) + \frac{c}{1-\gamma}}$

## Ex 3.16

For an episodic task, let the terminal timestep be $t = T$

Then,

$$\hat{G_t} = \sum_{k=0}^{T} \gamma^k \hat{R}_{t+k+1}$$

$$\hat{G_t} = \sum_{k=0}^{T} \gamma^k (R_{t+k+1} + c)$$

$$\hat{G_t} = \sum_{k=0}^{T} \gamma^k R_{t+k+1} + \sum_{k=0}^{T} c\gamma^k$$

$$\hat{G_t} = G_t + c\left(\frac{\gamma^{T+1} - 1}{\gamma - 1}\right)$$

$$\text{or} \quad \hat{G_t} = G_t + c\left(\frac{1 - \gamma^{T+1}}{1 - \gamma}\right) \quad \left[\text{since } \gamma < 1\right]$$

Following a similar approach for $\hat{v_\pi}(s)$

we get

$$\hat{v_\pi}(s) = E\left[ G_t + c\left(\frac{1 - \gamma^{T+1}}{1 - \gamma}\right) \,\middle|\, s_t = s \right]$$

Now, T is a random variable that depends on $s_t$

Thus, $c\left(\frac{1 - \gamma^{T+1}}{1 - \gamma}\right)$ cannot come out of the expectation

$\therefore$ a simple linear mapping doesn't exist between the two.

**Q5.** By Eqn 3.17

$$q_*(s,a) = E\left[R_{t+1} + \gamma V_*(s_{t+1}) \mid S_t = s, A_t = a\right]$$

Also, by eqn 3.18

$$V_*(s) = \max_a E\left[R_{t+1} + \gamma V_*(s_{t+1}) \mid S_t = s, A_t = a\right]$$

$$\therefore \boxed{V_*(s) = \max_a q_*(s,a)}$$

optimal value at state 's'

optimal value at state 's' after taking action 'a'