Q1)

Sol⁻ a) Firstly, we ~~can prove~~ need to prove

$$E_\pi\left[E_\pi\left[G_{t+1}|S_{t+1}\right]|S_t\right] = E_\pi\left[G_{t+1}|S_t\right]$$

Let $G_{t+1} = g'$, $S_{t+1} = s'$, $S_t = s$, $G_t = g$

$$\therefore E_\pi\left[E_\pi\left[g'|s'\right]|s\right]$$

$$= E_\pi\left[E_\pi\left[g'|s',s\right]\right]$$

$$= E_\pi\left[\sum_a \pi(a|s) \sum_{g'} g' \, P(g'|s',s)\right]$$

$$= \sum_{s'}\left[\sum_a \pi(a|s) \sum_{g'} g' \, P(g'|s',s)\right] \cdot P(s'|s)$$

$$= \sum_a \pi(a|s) \sum_{s'}\left[\sum_{g'} g' P(g'|s',s)\right] P(s'|s)$$

$$P(g'|s',s) = \frac{P(g',s',s)}{P(s',s)} \quad ; \quad P(s'|s) = \frac{P(s',s)}{P(s)}$$

$$\therefore E_\pi\left[E_\pi\left[g'|s'\right]|s\right] = \sum_a \pi(a|s) \sum_{s'} \sum_{g'} g' \cdot \frac{P(g',s',s)}{P(s',s)} \cdot \frac{P(s',s)}{P(s)}$$

$$= \sum_a \pi(a|s) \sum_{s'} \sum_{g'} (g' \, P(g'|s))$$

$$= E_\pi\left[g'|s\right]$$

$$= E_\pi\left[G_{t+1}|S_t\right] \quad —— Ⓐ$$

Now, we know,

$$V_\pi(s) = E_\pi\left[g_t + \gamma \sum_{k=0}^{\infty} \gamma^k R_{t+k+2} \Big| S_t = s\right]$$

$$= E_\pi\left[\eta + \gamma \, G_{t+1} | S_t = s\right]$$

$$= E_\pi\left[\eta|s\right] + \gamma E_\pi\left[G_{t+1}|S_t\right]$$

∵ from Ⓐ

$$V_n(s) = E_n[n|s] + \gamma E_n[E_n[G_{t+1}|S_{t+1}]|s_t]$$

$$= E_n[n + \gamma V_n(s')|s_t]$$

$$∴ V_n(s) = E_n[n + \gamma V_n(s')|s]$$

~~This is the state-value function of for Bellman~~

This is the Bellman's equations for state-value

b) For this, we need to prove firstly,

$$E_n[E_n[G_{t+1}|S_{t+1}, A_{t+1}]|S_t, A_t] = E_n[G_{t+1}|S_t, A_t]$$

Let $G_{t+1} = g'$ , $S_{t+1} = s'$ , $A_{t+1} = a'$, $G_t = g$, $S_t = s$, $A_t = a$

$$∴ E_n[E_n[g'|s',a']|s,a]$$

$$= E_n[E_n[g'|s',a',s,a]]$$

$$= E_n\left[\sum_a \pi \cdot (a|s) \sum_{g'} g' \, p(g'|s',a',s,a)\right]$$

$$= \sum_a \pi(a|s) \sum_{s/a'}\left[\sum_{g'} g' \, p(g'|s',a',s,a)\right] \cdot p(s',a'|s,a)$$

Ⓑ $$p(g'|s',a',s,a) = \frac{p(g',s',a',s,a)}{p(s',a',s,a)}$$

$$p(s',a'|s,a) = \frac{p(s',a',s,a)}{p(s,a)}$$

$$∴ E_n[E_n[g'|s'_a,a']|s,a] \oslash$$

$$= \sum_a \pi(a|s) \sum_{s',a'}\sum_{g'} g' \, \frac{p(g',s',a',s,a)}{p(s',a',s,a)} \cdot \frac{p(s',a',s,a)}{p(s,a)}$$

$$= \sum_a \eta(a|s) \sum_{g'} \sum_{s',a'} \frac{P(g',s',a'/s,a)}{P(s,a)}$$

$$= \sum_a \eta(a|s) \sum_{g'} P(g'|s,a)$$

$$= E_\pi \left[ g' | s,a \right]$$

$$= E_\pi \left[ G_{t+1} | S_t, a_t \right]$$

Now, we know, $q_\pi(s,a)$

$$= E_\pi \left[ \eta + \gamma \sum_{k=0}^{\infty} \gamma^k R_{t+k+2} \;\middle|\; S_t = s, A_t = a \right]$$

$$= E_\pi \left[ \eta + \gamma \, G_{t+1} | s,a \right]$$

$$= E_\pi \left[ \eta | s,a \right] + \gamma \, E_\pi \left[ G_{t+1} \, g' | s,a \right]$$

$$= E_\pi \left[ \eta | s,a \right] + \gamma \, E_\pi \left[ E_\pi \left[ g' | s',a' \right] | s',a \right]$$

$$= E_\pi \left[ \eta + \gamma \, q_\pi(s',a') | s,a \right]$$

$$\therefore q_\pi(s,a) = E_\pi \left[ \eta + \gamma \, q_\pi(s',a') \;\middle|\; S_t = s, A_t = a \right]$$

This is the Bellman's equation for action-value
of function.

Q2)

**61ᵗ a)** For the can collecting example, the state Transition table is as follows -

| S | S' | a | $p(s'|s,a)$ | $h(s,a,s')$ |
|------|------|---------|-------------|-------------|
| High | High | search | $\alpha$ | $r_{search}$ |
| High | Low | search | $1-\alpha$ | $r_{search}$ |
| High | High | wait | 1 | $r_{wait}$ |
| High | Low | wait | 0 | $r_{wait}$ |
| Low | Low | search | $\beta$ | $r_{search}$ |
| Low | High | search | $(1-\beta)$ | $-3$ |
| Low | Low | wait | 1 | $r_{wait}$ |
| Low | High | wait | 0 | $r_{wait}$ |
| Low | High | recharge | 1 | 0 |
| Low | Low | recharge | 0 | 0 |

For the other example, the state transition table is as follows -

| S | S' | a | P(s'|s, a) | $r(s,a,s')$ |
|---|----|---|-----------|-------------|
| Good | Good | Stay | ½ | + 3 |
| Good | Bad | Stay | ½ | − 1 |
| Good | Good | Move | 0 | + 3 |
| Good | Bad | Move | 1 | − 3 |
| Bad | Good | Stay | 0 | + 3 |
| Bad | Bad | Stay | 1 | − 1 |
| Bad | Good | Move | 1 | + 3 |
| Bad | Bad | Move | 0 | − 1 |

b) The state space diagram for can collecting bot example is as follows -



The state space diagram for the other example is as follows-