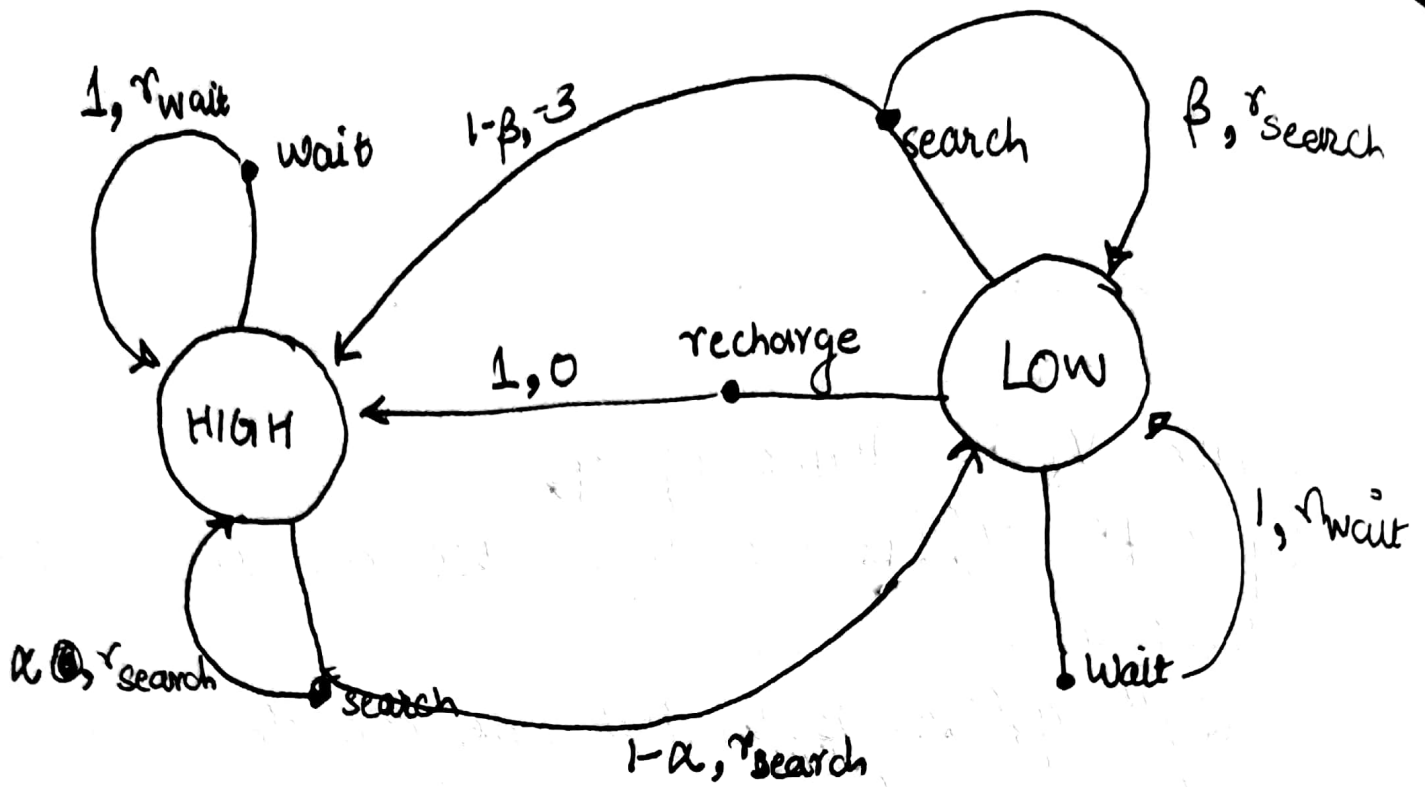


Q1



Let $P_{wait}(r)$ be a gaussian distribution with mean r_{wait} , standard deviation σ_{wait} .
 similarly $P_{search}(r) \sim N(r_{search}, \sigma_{search})$

S	a	S'	r	$P(S', r S, a)$
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	$1 * 1$

$$Q_3(a) \quad v_{\pi}(s) = E[G_t | S_t = s]$$

$$= \text{~~EV~~}$$

$$\text{Now, } G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots$$

$$= \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$$

Let \hat{r}_t be the new reward with ~~cost~~ constant c added to it.

$$\Rightarrow \hat{r}_t = r_t + c$$

$$\Rightarrow \hat{G}_t = \sum_{k=0}^{\infty} \gamma^k \hat{R}_{t+k+1}$$

$$= \sum_{k=0}^{\infty} [\gamma^k (R_{t+k+1}) + c \gamma^k]$$

$$= G_t + \sum_{k=0}^{\infty} c \gamma^k$$

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

$$\Rightarrow \hat{v}_{\pi}(s) = E[\hat{G}_t | S_t = s] = E[G_t + \frac{c}{1-\gamma} | S_t = s]$$

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

$$\Rightarrow \boxed{\hat{v}_{\pi}(s) = v_{\pi}(s) + \frac{c}{1-\gamma}}$$

Q(3)(b) If the task is episodic, let the task finish after $t=N$
 i.e. S_N is terminal state.

following from the above question

$$\begin{aligned}\hat{r}_t &= r_t + c \\ \hat{G}_t &= \sum_{k=0}^N \gamma^k \hat{R}_{t+k+1} = \sum_{k=0}^N \gamma^k (R_{t+k+1} + c) \\ &= \sum_{k=0}^N \gamma^k R_{t+k+1} + \sum_{k=0}^N c \gamma^k \\ &= G_t + c \left(\frac{\gamma^N - 1}{\gamma - 1} \right)\end{aligned}$$

$$\begin{aligned}\Rightarrow \hat{V}_\pi(s) &= E[\hat{G}_t | S_t = s] \\ &= \cancel{V_\pi(s) + c \left(\frac{\gamma^N - 1}{\gamma - 1} \right)} \\ &= E \left[G_t + c \left(\frac{\gamma^N - 1}{\gamma - 1} \right) | S_t = s \right]\end{aligned}$$

Since N is a random variable that depends on S_t , therefore it can't come out of expectation.

Hence, there is no simple mapping b/w $\hat{V}_\pi(s)$ and $V_\pi(s)$

Q5

write V_* in terms of q_*

by defⁿ $V_*^{(s)}$ is the best we can do starting at state s

by defⁿ $q_*(s, a)$ is the best we can do starting at state s and performing action a .

therefore,

$$V_*(s) = \max_{a \in A(s)} q_*(s, a)$$