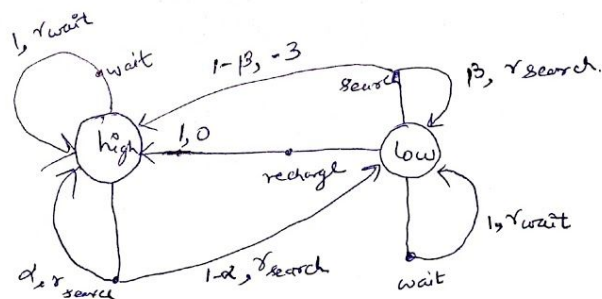


Assignment 2

Answer 1.

Table given below is.

s	a	s'	r	$P(s' \leq \sigma s, a)$
high	search	high	r_{search}	α
high	search	low	r_{search}	$1 - \alpha$
low	search	high	-3	$1 - \beta$
low	search	low	r_{search}	β
low	wait	low	r_{wait}	1
high	wait	high	r_{wait}	1
low	recharge	high	0	1



Answer 3, part 1

The value function of a state s under a policy π , denoted $V_{\pi}(s)$, is the expected return when starting in s and following π thereafter.

For MDPs V_{π} is defined as

$$\begin{aligned} V_{\pi}(s) &= E_{\pi} [G_t \mid S_t = s] \\ &= E_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s \right] \text{ for} \\ &\quad \text{all } s \in S \end{aligned}$$

Adding constant c to all the rewards

$$\hat{V}_{\pi}(s) = E_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k \hat{R}_{t+k+1} \mid S_t = s \right]$$

$$\text{where } \hat{R}_{t+k+1} = R_{t+k+1} + c$$

$$= E_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s \right]$$

$$+ E_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k c \mid S_t = s \right]$$

$$= V_{\pi}(s) + C \sum_{k=0}^{\infty} \gamma^k$$

$$= V_{\pi}(s) + \underbrace{C \times \frac{1}{1-\gamma}}_{\text{Constant} + V_c}$$

$$V_c = \frac{C}{1-\gamma}$$

Since V_c constant is added to the values of all states, and thus it does not affect the relative values of any state under any policies.

Answer 3 - part II

For episodic task. Let us assume task terminates after time step T .

$$\hat{v}_n(s) = v_n(s) + c \sum_{k=0}^T \gamma^k$$
$$= v_n(s) + c \left(\frac{1 - \gamma^{T+1}}{1 - \gamma} \right)$$

$$= v_n(s) + v_c$$

$$\text{where } v_c = c \left(\frac{1 - \gamma^{T+1}}{1 - \gamma} \right)$$

Since a constant is added to the values of all states, ~~it~~ like previous case it does not affect the relative values of any state under any policies.

Answer 5.

The optimal value functions are recursively related by the Bellman optimality equation.

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

Deterministic case

if,

$V_*(s) \leftarrow$ optimal value function

and $q_*(s, a) \leftarrow$ optimal action value function.

then

$$q_*(s, a) = r(s, a) + \gamma V_*(s')$$

$r(s, a) \leftarrow$ transition reward

$\gamma \leftarrow$ discounting factor.

$s' \leftarrow$ next state.

Stochastic case

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

when $p(s'|s, a)$ is the transition probability to new state s'

So,

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