

Name -
Ishaan Bassi

Assignment 2 Theory

Roll no - 2016238

Important →

where $p(s'|s,a) > 0$

i.e. $p(s'|s,a) = p(s|s,a)$

Q1. The table will stay the same except for the ~~few~~ rows where $p(s'|s,a) = 0$ are removed. This is due to the fact that only a single reward is possible with each triplet of (s, a, s') . Now we can see that for each (s, a) , $\sum p(s'|s,a) = 1$ as follows

high	S	a	
	high	search	$(\alpha + 1 - \alpha) = 1$
	low	search	$(1 - \beta + \beta) = 1$
	high	wait	$(1 + 0) = 1$
	low	wait	$(0 + 1) = 1$
	low	recharge	$(1 + 0) = 1$

Q6. (BUG)

The following two changes are made in the code -

- ① The policy ' π ' is taken to be ~~stoch~~ stochastic instead of deterministic i.e. all actions have some probability.
- ② If more than one action comes out to be optimal they are given equal probabilities.

Q 3. (a) Exercise 3.15

No, the signs of the rewards do not matter, instead only their values relative to each other are important.

This is due to the fact ~~for~~ that adding a constant to all the rewards does not change the relative distribution of state value function*. This can be seen as follows -

*Contd. In the continuing case.

Continuing Case

$$\begin{aligned} V_{\pi}(s_t) &= E[G_t | S_t = s] \\ &= E[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots | S_t = s] \end{aligned}$$

Adding a constant c to all rewards -

$$\begin{aligned} V'_{\pi}(s_t) &= E[R_{t+1} + c + \gamma(R_{t+2} + c) + \gamma^2(R_{t+3} + c) + \dots | S_t = s] \\ &= E[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots | S_t = s] \\ &\quad + c + \gamma c + \gamma^2 c + \dots \\ &= E[R_{t+1} + \gamma R_{t+2} + \dots | S_t = s] + \frac{c}{1-\gamma} \end{aligned}$$

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

\downarrow
 V_c

