



Q5. Give an equation for  $v^*$  in terms of  $q^*$ .

The optimal state-value function  $v^*$  is maximum over all policies,

$$v^*(s) = \max_{\pi} V_{\pi}(s)$$

where  $V_{\pi}(s) = E_{\pi}[G_t | S_t = s]$

The optimal action-value function  $q^*(s, a)$  is maximum over all  $q_{\pi}(s, a)$

$$q^*(s, a) = \max_{\pi} q_{\pi}(s, a)$$

where  $q_{\pi}(s, a) = E_{\pi}[G_t | S_t = s, A_t = a]$

$$q^*(s, a) = R_{t+1} + \gamma \sum_{s'} \pi(a|s') v^*(s')$$

$\downarrow$   
 $E[R_{t+1} | S_t = s, A_t = a]$

$$v^*(s) = \max_a q^*(s, a)$$

Q5. Exercise 3.15

$$v_{\pi}^*(s) = E_{\pi}[G_t | S_t = s]$$

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

Add a constant to  $R_{t+k+1}$ :  $\hat{R}_{t+k+1} = R_{t+k+1} + c$



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

$$= E_{\pi} [R_{t+1} + c + \gamma R_{t+2} + \gamma c + \gamma^2 R_{t+3} + \gamma^2 c + \dots \mid s_t = s]$$

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

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

$$= V_{\pi}(s) + c \underbrace{E \left[ \sum_{k=0}^{\infty} \gamma^k \right]}_{V_c}$$

$V_c$  is added to all states & hence relative value does not change.

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

Exercise 3.6 In episodic task, we have fixed no. of episodes, say  $k$ .

$$V_{\pi}(s) = E [R_{t+1} + \dots + \gamma^{k-1} R_{t+k} \mid s_t = s]$$

$$\bar{V}_{\pi}(s) = E [(R_{t+1} + c) + \gamma(R_{t+2} + c) + \dots + \gamma^{k-1} R_{t+k}] \mid s_t = s$$

$$\approx E[V_{\pi}(s)] +$$

$$= V_{\pi}(s) + c \left[ \frac{1 - \gamma^k}{1 - \gamma} \right]$$

Since, a constant value <sup>constant</sup> is added to all states, the relative importance of state value does not change.



Q1.

$$R_{t+1} = E [r_{t+1} | s_t = s, a_t = a, s_{t+1} = s']$$

↑ Expected reward at  $t+1$

$$= \sum_{r'} r' P\{r_{t+1}=r' | s_t=s, a_t=a, s_{t+1}=s'\}$$

$$P(r_{t+1}=r' | s_t=s, a_t=a, s_{t+1}=s') = \frac{P(s_{t+1}=s', r_{t+1}=r' | s_t=s, a_t=a)}{P(s_{t+1}=s' | s_t=s, a_t=a)}$$

$$P(s_{t+1}=s', r_{t+1}=r' | s_t=s, a_t=a) = P(r_{t+1} | s_t, a_t, s_{t+1}) \times P(s_{t+1} | s_t, a_t)$$

We are given  $P(r | s, a, s')$  &  $P(s' | s, a)$  in the table.

$s$	$a$	$s'$	$p(s'   s, a)$	$r(s, a, s')$	$p(s', r   s, a)$
high	search	high	$\alpha$	$r_{\text{search}}$	$\alpha r_{\text{search}}$
high	search	low	$1-\alpha$	$r_{\text{search}}$	$(1-\alpha) r_{\text{search}}$
low	search	high	$1-\beta$	$-3$	$-3(1-\beta)$
low	search	low	$\beta$	$r_{\text{search}}$	$\beta r_{\text{search}}$
high	wait	high	1	$r_{\text{wait}}$	$r_{\text{wait}}$
high	wait	low	0	-	0
low	wait	high	0	-	0
low	wait	low	1	$r_{\text{wait}}$	$r_{\text{wait}}$
low	recharge	high	1	0	0
low	recharge	low	0	-	0