

Ques 1.

APoorv KHATTAR  
2016016

s	a	s'	$\pi$	$P(s', \pi   s, a)$
high	search	high	$\pi_{\text{search}}$	$\alpha$
high	search	low	$\pi_{\text{search}}$	$1-\alpha$
high	wait	high	$\pi_{\text{wait}}$	1
low	search	low	$\pi_{\text{search}}$	$\beta$
low	search	<del>low</del> /high	<del><math>\pi_{\text{search}}</math></del> (-3)	$(1-\beta)$
low	wait	low	$\pi_{\text{wait}}$	1
low	recharge	high	0	1

←  
Exercise 3.4

Ques 5.  $V_{\pi}(s) = \max_{a \in A(s)} q_{\pi}(a, s)$

$$\begin{aligned}
 &= \max_a E_{\pi}[G_t | A_t=a, S_t=s] \\
 &= \max_a E_{\pi}[R_t + \gamma G_{t+1} | A_t=a, S_t=s] \\
 &= \max_a E_{\pi}[R_t + \gamma V_{\pi}(S_{t+1}) | A_t=a, S_t=s] \\
 &= \max_a \sum_{s'} \sum_{\pi} P(s', \pi | s, a) \{ \pi + \gamma V_{\pi}(s') \} \\
 &= \max_a \sum_{s'} \sum_{\pi} P(s', \pi | s, a) \{ \pi + \gamma \max_a q_{\pi}(s', a) \}
 \end{aligned}$$

Ques 3.

(a) Exercise 3.15

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \quad (1)$$

If we add constant to all rewards then,

$$\begin{aligned}
 G'_t &= (R_{t+1} + c) + \gamma (R_{t+2} + c) \dots = \sum_{k=0}^{\infty} \gamma^k (R_{t+k+1} + c) \\
 &= \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} + \sum_{k=0}^{\infty} \gamma^k c
 \end{aligned}$$

$$\Rightarrow G'_t = G_t + c \left( \frac{1}{1-\gamma} \right)$$

$$\Rightarrow G'_t = G_t + \gamma c$$

Let  $V_{\pi'}(s)$  be the new value function for  $G_{t'}$  then

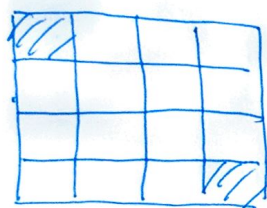
$$\begin{aligned} V_{\pi'}(s) &= E_{\pi'} [G_{t'} | S_t = s] \\ &= E_{\pi'} [G_t + V_c | S_t = s] \\ &= E_{\pi'} [G_t | S_t = s] + E_{\pi'} [V_c | S_t = s] \\ &= V_{\pi}(s) + V_c \end{aligned}$$

$\Rightarrow V_{\pi'}(s)$  increases by  $V_c$  and  $V_c = \frac{c}{1-\gamma}$

(b) Exercise 3.16

Adding a constant ( $c$ ) to all rewards in an episodic task will effect the tasks.

(considers the following episodic task:-



(0,0) and (3,3) are terminal states  
reward for each action (-1)

Optimal Policy would ensure shortest path to either of the terminal states.

Now if the rewards are  $-1+c$  ( $c > 1$ ) then optimal policy changes, also value functions will for any policy will be different as it would never try to go to terminal states for maximizing reward