

Q1

State set = {high, low}

$A(\text{high}) = \{\text{search, mail}\}$

$A(\text{low}) = \{\text{search, mail, recharge}\}$

$s$	$a$	$s'$	$r$	$p(s'   s, a)$
high	search	high	$r_{\text{search}}$	$\alpha$
high	search	low	$r_{\text{search}}$	$1 - \alpha$
high	mail	high	<del><math>r_{\text{search}}</math></del> $r_{\text{wait}}$	1
high	mail	low	<del><math>r_{\text{search}}</math></del> <del><math>r_{\text{wait}}</math></del>	—
low	search	high	-3	$1 - \beta$
low	search	low	$r_{\text{search}}$	$\beta$
low	mail	low	$r_{\text{mail}}$	1
low	mail	high	—	—
low	recharge	high	0	1
low	recharge	low	—	—

Q3

Rewards are +ve ~~reward~~ for goals, -ve for running into the edge of the world and zero the rest of the time

Only the intervals between the rewards are important and not their signs.

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

$$= \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \quad \gamma = \text{discount rate}$$

$$0 \leq \gamma \leq 1$$

Using the above equation, one can show that adding a constant  $c$  to all the rewards adds a constant  $V_c$  to the values of all the states and hence does not affect the relative values of any states.

Here rewards can be made the same sign by adding/subtracting a large positive quantity from to/from all the rewards.

which leads to an increase or ~~an~~ a decrease in the value function by a constant. This whole thing ~~is~~ has no effect on the algorithm.

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

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

$$= E\left[\left(\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} + \sum_{k=0}^{\infty} \gamma^k c\right) \mid S_t = s\right]$$

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

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

$$= V_{\pi}(s) + \frac{c}{1-\gamma} \quad \left. \vphantom{\frac{c}{1-\gamma}} \right\} \text{constant term } V_c$$



Q5

Equation for  $V_*$  in terms of  $q_*$

$$\begin{aligned} q_*(s, a) &= E \left[ R_{t+1} + \gamma V_*(S_{t+1}) \mid S_t = s, A_t = a \right] \\ &= \sum_{s', r} p(s', r \mid s, a) \left[ r + \gamma \max_{a'} q_*(s', a') \right] \end{aligned}$$

From the Bellman Optimality equation

$$\begin{aligned} V_*(s) &= \max_{a \in A(s)} \pi_*(s, a) \\ &= \max_a \sum_{s', r} p(s', r \mid s, a) \cdot [r + \gamma V_*(s')] \\ &= \max_a \sum_{s', r} p(s', r \mid s, a) \left[ r + \gamma \max_{a'} q_*(s', a') \right] \end{aligned}$$

Q3:16(b)

In case of episodic task

let terminal time be  $T$

solving the equation

$$\begin{aligned} V_{\pi}(s) &= E \left[ G_t' \mid S_t = s \right] \\ &= E \left[ \sum_{k=0}^T \gamma^k (R_{t+k+1} + c) \mid S_t = s \right] \\ &= E \left[ \sum_{k=0}^T \gamma^k R_{t+k+1} \mid S_t = s \right] + \\ &\quad E \left[ \sum_{k=0}^T \gamma^k c \mid S_t = s \right] \\ &= V_{\pi}(s) + E \left[ \sum_{k=0}^T \gamma^k c \mid S_t = s \right] \end{aligned}$$

$V_c = E \left[ \sum_{k=0}^T \gamma^k c \mid S_t = s \right]$  is a function of  $T$

$\therefore T$  is a random variable, which normally varies from episode to episode.

$\Rightarrow$  diff episodes will have diff value functions

$$G_t \text{ is } G_t + c \left( \frac{1 - \gamma^T}{1 - \gamma} \right)$$

Thus it will increase  $V_{\pi}$  when  $T$  will increase