

## Assignment 2

Neelam R.R.  
PHD 19108.

① States = {High, low}

Actions = {search, wait, Recharge}

$$p(s' = \text{high} | s = \text{high}, a = \text{search}) = \alpha$$

$$p(s' = \text{low} | s = \text{high}, a = \text{search}) = 1 - \alpha$$

$$p(s' = \text{high} | s = \text{low}, a = \text{search}) = 1 - \beta.$$

⋮

The table can be filled as.

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

② Exercise 3.15

③ Eqn (3.8)

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

Where  $\gamma$  is a parameter,  $0 \leq \gamma \leq 1$  called the discount rate.

$$G_t = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$$

$$G_t^c = \sum_{k=0}^{\infty} \gamma^k (R_{t+k+1} + C_{t+k+1})$$

$$= \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} + C_{t+k+1} \sum_{k=0}^{\infty} \gamma^k$$

$$C_{t+k+1} = C$$

$$= \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} + \frac{C}{1-\gamma}$$

$$= \underline{\underline{G_t + \frac{C}{1-\gamma}}}$$

$$V_c = \mathbb{E}[G_t^c | s_t = s]$$

$$= \mathbb{E}\left[G_t + \frac{C}{1-\gamma} \mid s_t = s\right]$$

$$= \mathbb{E}[G_t | s_t = s] + \mathbb{E}\left[\frac{C}{1-\gamma}\right]$$

$$= \underline{\underline{\mathbb{E}[G_t | s_t = s] + \frac{C}{1-\gamma}}}$$

$$\mathbb{E}[\text{constant}]$$

= constant

$\therefore V^c$  does not affect the relative values of any states.

### Exercise 3.16.

$$\text{Episodic task} :- G_t = \sum_{k=t+1}^T \gamma^{k-t-1} R_k.$$

The rewards, we will sum over all the rewards to find the best action to be picked.

When adding a constant,

$$G^c = \sum_{k=t+1}^T \gamma^{k-t-1} (R_k + c)$$

The rewards obtained in the episodic task is very important. Adding a constant can affect the rewards obtained in a state, and affect the action to be picked in the state.

eg: If  $R = -5$  &  $c = 2$ .

$$\begin{aligned} \text{The reward at a state } s &= -5 + 2 \\ &= \underline{\underline{-3}} \end{aligned}$$

The reward value increased, and chances of picking a non-best arm is high.

$$\begin{aligned} q_{\pi}(s, a) &= \mathbb{E}[G_t | s_t = s, A_t = a] \\ &= \mathbb{E}\left[\sum_{k=t+1}^T \gamma^{k-t-1} (R_k + c) \mid s_t = s, A_t = a\right] \end{aligned}$$

$\therefore$  The action value function will also change.

## ⑤ Optimal state value function $V^*$

A policy  $\pi$  is defined to be better than or equal to a policy  $\pi'$  if its expected return is greater than or equal to that of  $\pi'$  for all states.  $\pi \geq \pi'$  if and only if  $V_\pi(s) \geq V_{\pi'}(s)$ . The optimal policy is that policy which is better than or equal to all other policies. It share a ~~same~~ same state value function denoted by  $V^*$ .

$$V^*(s) = \max_{\pi} V_{\pi}(s) \quad \forall s \in S.$$

Optimal action value function  $q^*$  gives the expected return for taking action  $a$  in state  $s$  and thereafter following an optimal policy.

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

$$\boxed{V^*(s) = \max_{a \in A(s)} q_{\pi^*}(s, a)}$$

It is taking the best action by taking best policy  $\pi^*$  into account and giving the best expected return.



⑧ Yes, we can express  $R_{t+2}$  in terms of  $S_t$  &  $A_t$ . That means  $R_{t+2}$  is dependent on  $S_t$  &  $A_t$ . This happens when the previous state of  $S_{t+2}$  i.e.  $S_{t+1}$  doesn't produce any reward or didn't take any action.

$$\therefore p(S_{t+2}=s'', R_{t+2}=r'' | S_t=s, A_t=a) =$$

$$Pr\{S_{t+2}=s'', R_{t+2}=r'' | S_t=s, A_t=a\}$$

$$v(s, a, s'') = \mathbb{E}[R_{t+2} | S_t=s, A_t=a, S_{t+2}=s'']$$

$$= \sum_{r'' \in R} r'' \frac{p(s'', r'' | s, a)}{p(s'' | s, a)}$$

=====

This happens when the previous state of  $R_{t+2}$  i.e.  $S_{t+1}$  does not happen.

$$\textcircled{9} \mathbb{E}[R_{t+2} | S_t=s, A_t=a] = \mathbb{E}_{s''} [\mathbb{E}[R_{t+2} | S_{t+2}=s'', s, a] | s, a]$$

$$= \mathbb{E}_{s''} \left[ \sum_{r''} r'' \frac{p(s'', r'' | s, a)}{p(s'' | s, a)} \middle| s, a \right]$$

$$= \mathbb{E}_{s''} [r''(s, a, s'') | s, a]$$

$$s'' = s_{t+2}$$

$$r'' = R_{t+2}$$

$$= \sum_{s'' \in \mathcal{S}} \cancel{P''} r''(s, a, s'')$$

=

$$\therefore \mathbb{E}[R_{t+2} | s_t = s, A_t = a] = \mathbb{E}_{s_{t+2}''} [r''(s, a, s'') | s, a]$$

⑩ The state value function

$$V_{\pi}(s) = \mathbb{E}_{\pi} [G_t | s_t = s]$$

As we know,  $G_t = R_{t+1} + \underbrace{\gamma G_{t+1}}_{\text{Reward to go function.}}$

If  $\gamma = 1$ , the agent is far-sighted of the future.

$$\therefore V_{\pi}(s) = \mathbb{E}_{\pi} [R_{t+1} + \gamma G_{t+1} | s_t = s]$$

$$= \sum_a \mathbb{E} [R_{t+1} | s_t = s] + \gamma \mathbb{E}_{\pi} [G_{t+1} | s_{t+1} = s']$$

$$= \sum_{s'} \sum_a \gamma + \gamma \mathbb{E}_{\pi} [G_{t+1} | s_{t+1} = s']$$

Summing over all actions, to get the state value function.

$$= \sum_a \pi(a|s) \sum_{s', r} p(s', r|s, a) [r + \gamma V \pi(s')] \quad \forall s \in S.$$

This is the sum of all values of three variables  $a, s', r$ .

For each value of  $a, s', r$ , we compute the probability, i.e.  ~~$\pi(a)$~~   $\pi(a|s) p(s', r|s, a)$  & multiply with the current reward & reward to go term.

Then we sum over all possibilities to get an expected value.

② An agent receives a sequence of rewards  
 $R_1 = 2, R_2 = -1, R_3 = 10$  &  $R_4 = -3$ .

$$\gamma = 0.5$$

The infinite horizon discounted return is

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

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

$$G_t = R_{t+1} + \gamma G_{t+1}$$

$$G_4 = 0.$$

$$G_3 = R_{t+4} = R_4 = -3.$$

$$G_2 = R_3 + \gamma G_3 = -3 \times 0.5 + 10 \\ = 8.5$$

$$G_1 = R_2 + \gamma G_2 = -1 + 0.5 \times 8.5 \\ = \underline{\underline{3.25}}$$

$$G_0 = R_1 + \gamma G_2 = 2 + 0.5 \times 3.25 \\ = \underline{\underline{3.625}}.$$

The  $\gamma$  discounted reward/return for each time step.

$$G_0 = 3.625$$

$$G_1 = 3.25$$

$$G_2 = 8.5$$

$$G_3 = -3$$

$$G_4 = 0.$$



(12) We know, the optimal state value function

$$v^*(s) = \mathbb{E}_{\pi^*} [R_{t+1} + \gamma v^*(s_{t+1}) | s_t = s]$$

By Bellman's optimality principle,

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

As we have given  $v^*(s)$ ,

$$v^*(s) = \max_a q^*(s, a) = \max_a [R_{t+1} + \gamma v^*(s_{t+1}) | s_t = s, A_t = a]$$

$$\therefore \pi^*(s) = \operatorname{argmax}_{a \in A(s)} \mathbb{E} [R_{t+1} + \gamma v^*(s_{t+1}) | s_t = s, A_t = a]$$

$$\pi^*(s) = \operatorname{argmax}_{a \in A(s)} \sum_{s', r} p(s', r | s, a) [r + v^*(s')]$$

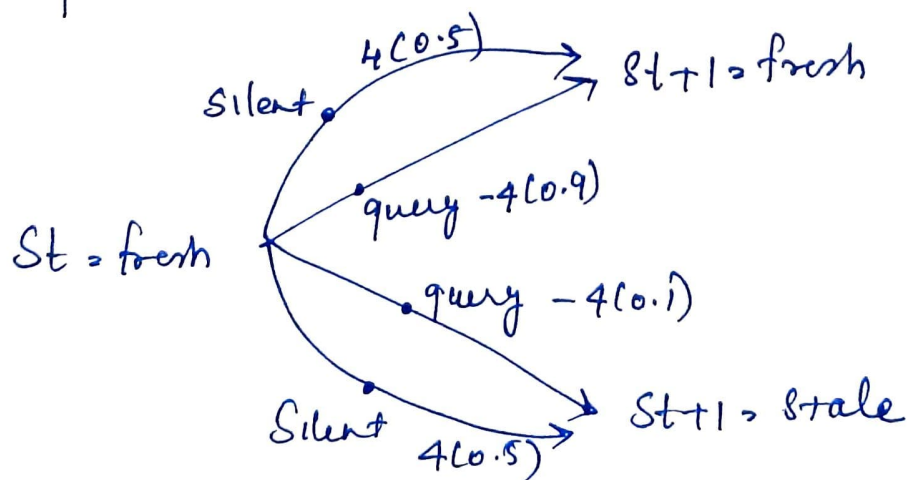
---

13

(i) MDP.

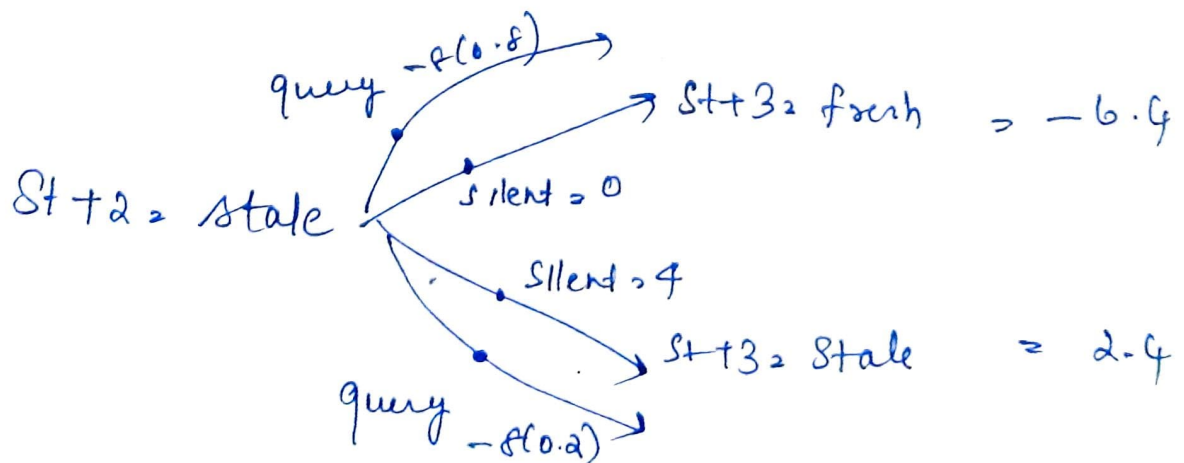
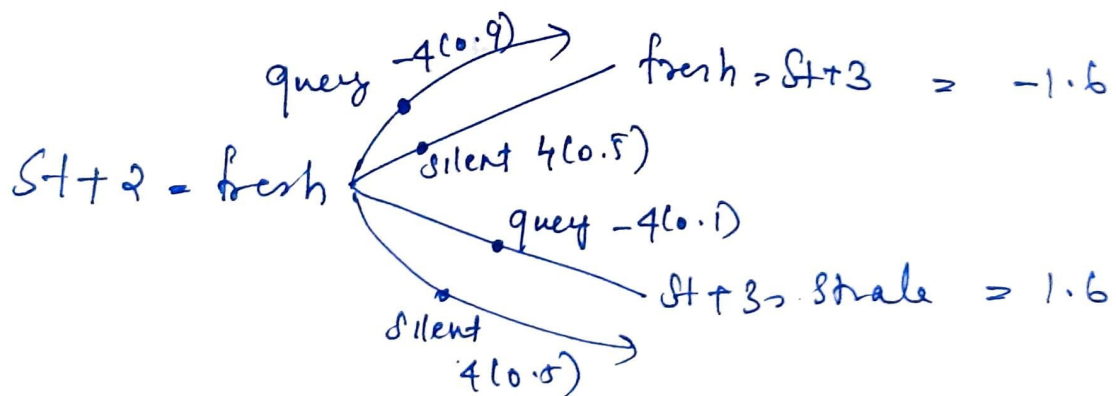
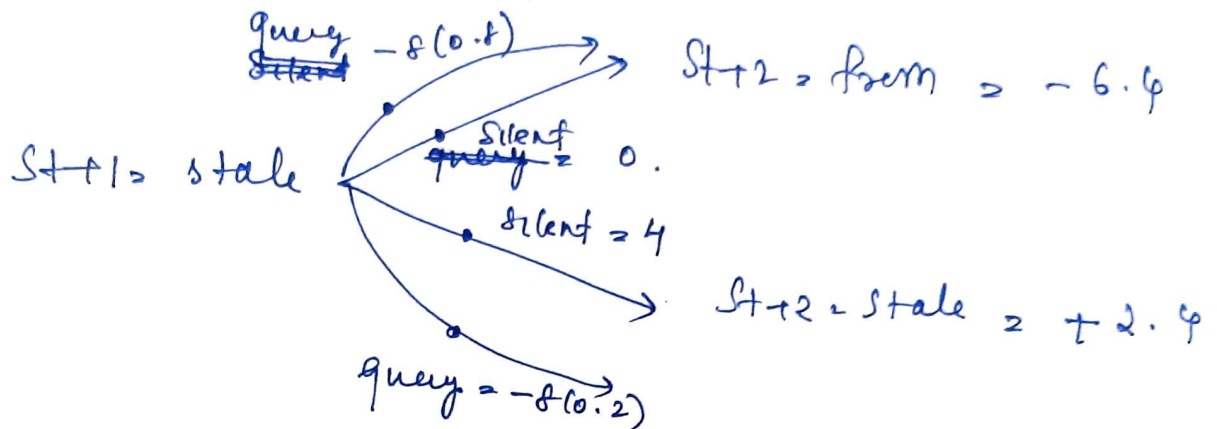
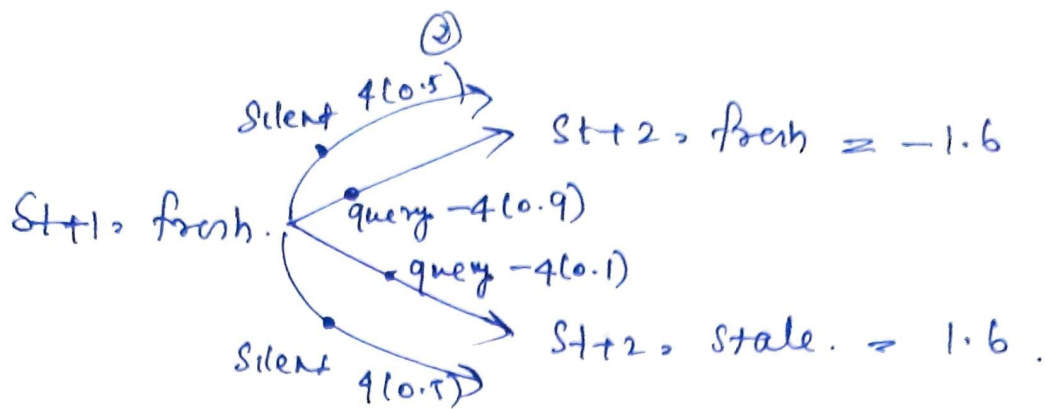
$S_t$	$A_t$	$S_{t+1}$	$r_{t+1}$	$P$
fresh	query	fresh	-4	0.9
fresh	query	stale	-4	0.1
stale	query	fresh	-8	0.8
stale	query	stale	-8	0.2
fresh	silent	fresh	+4	0.5
fresh	silent	stale	+4	0.5
stale	silent	<del>stale</del> fresh	+4	0
stale	silent	stale	+4	1

dis Starting with state = fresh.



$$S_{t+1} = \text{fresh} = -1.6.$$

$$S_{t+1} = \text{stale} = 1.6.$$



(3)

 $S_t \quad S_{t+1} \quad S_{t+2} \quad S_{t+3}$ 
 $\text{Fresh}(F) \text{ Fresh}(F) \text{ Fresh}(F) \text{ Fresh}(F)$ 
 $F \quad F \quad F \quad \text{Stale}(S)$ 
 $F \quad F \quad S \quad S$ 
 $F \quad F \quad S \quad F$ 
 $F \quad S \quad F \quad F$ 
 $F \quad S \quad S \quad F$ 
 $F \quad S \quad F \quad S$ 
 $F \quad S \quad S \quad S$ 

$$R_t = R_{t+1} + \gamma G_{t+1}$$

$$-1.6 + \frac{1}{2}(-1.6 - 1.6) + 10$$

$$= \underline{\underline{6.8}}$$

$$-1.6 + \frac{1}{2}(-1.6 + 1.6) + 10$$

$$= \underline{\underline{8.4}}$$

$$-1.6 + \frac{1}{2}(+1.6 + 2.4) + 10$$

$$= \underline{\underline{10.4}}$$

$$-1.6 + \frac{1}{2}(1.6 + 6.4) + 10$$

$$= \underline{\underline{6}}$$

$$1.6 + \frac{1}{2}(-6.4 - 1.6) + 10$$

$$= \underline{\underline{7.6}}$$

$$1.6 + \frac{1}{2}(2.4 + 6.4) + 10$$

$$= \underline{\underline{9.6}}$$

$$1.6 + \frac{1}{2}(-6.4 + 1.6) + 10$$

$$= \underline{\underline{9.2}}$$

$$1.6 + \frac{1}{2}(2.4 + 2.4) + 10$$

$$= \underline{\underline{-6}}$$

The optimal policy is that giving maximum reward

$$\therefore \pi^* =$$

$$S_t = \text{Fresh}$$

$$S_{t+1} = \text{Fresh}$$

$$S_{t+2} = \text{Stale}$$

$$S_{t+3} = \text{Stale}.$$

(14) Policy improvement step:-

- (i) Either improves policy or leaves it unchanged.
- (ii) If it leaves unchanged, the policy is optimal policy.

The operator  $T_{\pi} f(s) = \mathbb{E}[R_{t+1} + \gamma f(s_{t+1}) | s_t]$

$$T_{\pi} V_{\pi_k}(s) = V_{\pi_k}(s) \quad s \in S.$$

$$\therefore T[V_{\pi_k}] = T_{\pi_{k+1}} V_{\pi_k}.$$

$$\text{i.e. } V_{\pi_{k+1}}(s) \geq V_{\pi_k}(s)$$

Proof:-

$$T_{\pi_{k+1}} V_{\pi_k}(s) = T V_{\pi_k}(s) \geq T_{\pi_k} V_{\pi_k}(s) = V_{\pi_k}(s)$$

$$T_{\pi_{k+1}} V_{\pi_k}(s) \geq V_{\pi_k}(s)$$

$$\therefore T_{\pi_{k+1}}(T_{\pi_{k+1}}(V_{\pi_k}(s))) \geq V_{\pi_k}(s)$$

$$\therefore T_{\pi_{k+1}}^N V_{\pi_k}(s) \geq V_{\pi_k}(s) \quad \text{--- (1)}$$

Applying limit

$$\lim_{N \rightarrow \infty} T_{\pi_{k+1}}^N V_{\pi_k}(s) = V_{\pi_{k+1}}(s) \quad \text{--- (2)}$$

From (1), (2) can be written as

$$\lim_{N \rightarrow \infty} T_{\pi_{k+1}}^N V_{\pi_k}(s) \geq V_{\pi_k}(s)$$



$$\therefore \boxed{V_{\pi_{k+1}}(s) \geq V_{\pi_k}(s)} \quad \forall s \in S$$

$$\nexists V_{\pi_{k+1}}(s) = V_{\pi_k}(s).$$

$\therefore T_{\pi_{k+1}}$  has a unique fixed point  $V_{\pi_{k+1}}$

$$\therefore \boxed{V_{\pi_k} = V_{\pi_{k+1}}} \quad \forall s \in S.$$

(15)

$n = 2$

(1) Iteration

2nd iteration

	$s_t$	$a_t$	$s_{t+1}$	$r_{t+1}$	$P$	$s_{t+2}$	$r_{t+2}$	$P$
Step 1	up		2 (up)	-1	0.5	F	4	0.5
Step 1	down		G	1	0.5	-	-	-
Step 2	up		F	4	0.5	-	-	-
Step 2	down		1 (up)	2	0.5	G	2	0.5

(1) Iteration  $\Rightarrow$  Reward  $\rightarrow R_1$

$$(-1 \times 0.5 + 1 \times 0.5) + 1(4 \times 0.5) \\ = \underline{\underline{2}}.$$

(2) Iteration  $\Rightarrow$  Reward  $R_2$ .

$$(4 \times 0.5 + 2 \times 0.5) + 1(2 \times 0.5) \\ = \underline{\underline{4}}.$$

So as we are increasing the iterations, or increasing the value of  $n$ , the iterations increase proportional. If we are considering 'n' steps, there will be 'n' iterations.

The optimal policy is given by.

Step 2  $\longrightarrow$  Step 1  $\longrightarrow$  Ground.