

Name : SAHAS MARWAH

Roll No. : 2020237

RL HW -2

1. $p(\text{even}) = 2 p(\text{odd})$, Gaussian ($i, 1$)

1. $\rightarrow 10 \text{ items / arms} \rightarrow 10 \text{ times}$

$$E[\text{total reward after 10 pulls}] = 10 \times E[q_*(a)]$$

and, $E[q_*(a)]$

$$= \sum p(a) \cdot q_*(a)$$

$$= \frac{2}{15} (2+4+6+8+10) + \frac{1}{15} (1+3+7+5+9)$$

$$= 4 + \frac{5}{3} = \frac{17}{3}$$

$$\left\{ \begin{array}{l} p(a) = \\ \quad \begin{cases} 2/15, & a = \text{even} \\ 1/15, & a = \text{odd} \end{cases} \\ b/w [1, 10] \end{array} \right.$$

$\in I$.

Hence, $10 E[q_*(a)]$

$$= \frac{170}{3} = 56.66$$

2. $R=0 \mid p=1/2 \quad \forall a \in \{1, 2, 4, 5, 7, 9, 10\}$

$$R=1 \mid p \neq 1/2$$

$$R=0 \mid p = 3/10$$

$$R=0.2 \mid p = 3/10$$

$$R=1 \mid p = 4/10$$

Now, $q_*^{(1)}(a) = 0(0.5) + 1(0.5) = \underline{0.5}$

and, $q_*^{(2)}(a) = 0(0.3) + 0.2(0.3) + 1(0.4)$
 $= \underline{0.46}$

Now, for optimal stochastic policies we will make 6 combinations of arms with $q_*(a) = 0.5$.

1. Arms (1, 2, 4, 5, 7) with ~~0.2~~ 0.2 each
 2. Arms (1, 2, 4, 5) with ~~0.25~~ 0.25 each
 3. Arms (1, 2, 4) with ~~1/3~~ $\frac{1}{3}$ each
 4. Arms (1, 2) with ~~1/2~~ $\frac{1}{2}$ each
 5. Arms (1, 2, 4, 5, 7, 9) with ~~1/6~~ $\frac{1}{6}$ each
 6. Arms (1) with 0.4 and Arm(2) with 0.6.
-
4. Table for $p(s', r|s, a)$

s	a	s'	r'	$p(s', r s, a)$
high search	high recharge	high search	-3	α
high search	low recharge	high search	1 - α	$(1 - \alpha)$
low search	high recharge	-3	1 - β	$(1 - \beta)$
low search	low recharge	high search	β	β
high wait	high recharge	high wait	1 - γ	$1 - \gamma$
low wait	low recharge	high wait	γ	γ
low recharge	high wait	0	1	1

* This table was derived by MDP table given in example 3.3. and the transition graph.

* As rewards have no dependence on prob. distribution we can say that $p(s', r|s, a) = p(s'|s, a)$.

3. $\gamma \in \{0, 1\}$. explore at $t=1, 3, 5$ (over non-greedy only)
 $Q_1(a), a \in \{1, 2, 3\}$ exploit at $t=2, 4, 6$

Sample Mean: let $Q_1(1) = 0.2$
 $Q_1(2) = 0.3$

At $t=1$: $Q_1(3) = 0.5$

let $A_1 = 1, R_1 = 1$

$Q_2(1) = \frac{0.2 + 1}{2} = 0.6, Q_2(2) = 0.3, Q_2(3) = 0.5$

At $t=2$:

(Greedy) let $A_2 = 1, R_2 = 0$

$Q_3(1) = \frac{0.6 + 0}{3} = 0.2, Q_3(2) = 0.3, Q_3(3) = 0.5$

At $t=3$:

let $A_3 = 2, R_3 = 0$

$Q_4(1) = 0.2, Q_4(2) = \frac{0.3 + 0}{2} = 0.15, Q_4(3) = 0.5$

At $t=4$:

(Greedy) let $A_4 = 3, R_4 = 1$

$Q_5(1) = 0.2, Q_5(2) = 0.15, Q_5(3) = \frac{0.5 + 1}{2} = 0.75$

At $t=5$:

let $A_5 = 1, R_5 = 1$

$Q_6(1) = \frac{0.2 + 1}{4} = 0.3, Q_6(2) = 0.15, Q_6(3) = 0.75$

(Greedy) At $t=6$:

let $A_6 = 3, R_6 = 0$

$Q_7(1) = 0.3, Q_7(2) = 0.15, Q_7(3) = \frac{0.75 + 0}{3} = 0.25$

6. 3.15 show that intervals between rewards do not matter.

The signs of rewards is not important, only the intervals between them are as the difference of rewards matters.

Generally, if good reward = $-1, 2, 5$
and bad reward corresponding = $-3, 3$
They will ~~be~~ have same results.

Now,

$$\text{Eq. 3.8} \Rightarrow G_t = \sum_{k=0}^{\infty} \gamma^k R_{t+k} + \dots$$

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

Adding a const. c to all rewards:

$$\text{using } G_t = \sum_{k=0}^{\infty} \gamma^k [R_{t+k} + c]$$

$$\therefore G_t = \frac{1}{1-\gamma} (R_{t+1} + c)$$

As every state reward gets an additional const. term, added to it, $v_c = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k c \right]$

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

Hence, it does not affect relative values in any state as $v_c = \frac{c}{1-\gamma}$ in terms of c and γ .

3.16

In an episodic task, adding a const. to all rewards does affect the agent because the cumulative reward is dep. on the length of the episode.

Let us say if $r = -1$, the length of the episode decreases as the agent wants to find the exit quickly (in maze runner).

But if say $r = +1$ because of some +ve value of c , the length of episode increases as agent does not want to find the exit in maze runner.

* This case will go onto ∞ .

$$s' = p(s', r|s, a)$$

$$(s', a) \sim \pi(a'|s) \cdot p(s', r|s, a)$$

classmate

Date _____
Page _____

8. $v_\star(s)$ in terms of q_\star :

$$v_\star(s) = \max_{a \in A(s)} q_\star(s, a)$$

$$= \max_{a \in A(s)} \sum_{s', r} p(s', r|s, a) [r + \gamma \max_{a' \in A(s')} q_\star(s', a')]$$

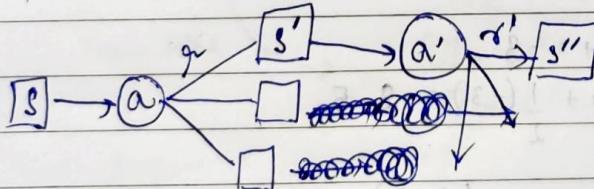
11. RV R_{t+2} dependent on $| s_t, A_t = ?$

$$p(R_{t+2} | s', a') = \sum_{s''} p(R_{t+2}, s'' | s', a')$$

- Now, $p(R_{t+2} | s, a) = \sum_{s', a'} p(s', a' | s, a) \cdot p(R_{t+2} | s', a')$

$$= \sum_{s', a'} [\pi(a'|s') \cdot p(s', r|s, a)] \cdot \sum_{s''} [p(R_{t+2}, s'' | s', a)]$$

$$= \sum_{s', r} p(s', r | s, a) \cdot \sum_{s'', a'} \pi(a'|s') p(s'', r' | s', a)$$



12. $E[R_{t+2} | s_t = s, A_t = a]$

$$= \sum_{r'} r' P(r = R_{t+2} | s_t = s, A_t = a)$$

(from above)

$$= \sum_{r'} r' \sum_{s', a} p(s', r | s, a) \sum_{s'', a'} \pi(a' | s') p(s'', r' | s', a)$$

↳ MDP's PMF

$$\begin{aligned}
 13. \quad v_{\pi}(s) &= E[G_t | s_t = s] \\
 &= E_{\pi}[R_{t+1} + \gamma G_{t+1} | s_t = s] \\
 &= E_{\pi}[R_{t+1} | s_t = s] + \gamma E_{\pi}[G_{t+1} | s_t = s] \\
 &= E_{\pi}[R_{t+1} | s_t = s] + \gamma E_{\pi}[R_{t+2} + \gamma G_{t+2} | s] \quad \text{recursively} \\
 &= E_{\pi}[R_{t+1} | s] + \gamma E_{\pi}\left[\sum_a \pi(a | s) \sum_{s'} p(s' | s, a) [r + \gamma E_{\pi}[G_{t+1} | s_{t+1} = s']] \right] \\
 &= \sum_a \pi(a | s) \sum_{s', r} p(s' | s, a) [r + \gamma v_{\pi}(s')] \\
 &\quad \left(\text{Bellman Equation for } v_{\pi}(s) \right)
 \end{aligned}$$

$$14. \quad R_1 = 2, R_2 = -1, R_3 = 10, R_4 = -3, \gamma = 0.5$$

Discounted reward for each timestep:

$$G_0 = R_1 + \gamma G_1$$

$$G_1 = R_2 + \gamma G_2$$

$$G_2 = R_3 + \gamma G_3$$

$$G_3 = R_4 + \gamma G_4 \quad \therefore G_3 = R_4 = -3$$

$$\Rightarrow G_2 = 10 + \frac{1}{2}(-3) = 8.5$$

$$\Rightarrow G_1 = -1 + \frac{1}{2}(8.5) = 3.25$$

$$\Rightarrow G_0 = 2 + \frac{1}{2}(3.25) = 3.625$$

$$\begin{aligned}
 \text{New } G_t &= R_t + \gamma R_{t+1} + \gamma^2 R_{t+2} + \dots + \infty \\
 &= c + \gamma c + \gamma^2 c + \dots + \infty
 \end{aligned}$$

$$(c + \gamma c + \gamma^2 c) = c \left(\frac{1}{1-\gamma} \right) = \frac{c}{1-\gamma}$$

15. Given: $v_{\pi}(s) \neq \text{sgn}$

We know, $\pi_{\star}(s) = \max_a q_{\pi}(s, a)$

$$\pi_{\star}(s) = \max_a \sum_{s', r} p(s', r | s, a) [r + \gamma v_{\pi}(s')]$$

so, in this case we take the best action that returns the highest value to maintain optimality by iterating over all states actions that can be chosen over a state greedily.

17. Given: $v_{\pi'}(s) \geq v_{\pi}(s)$

Prove : $v_{\pi'}(s) \geq v_{\pi}(s)$

We know,

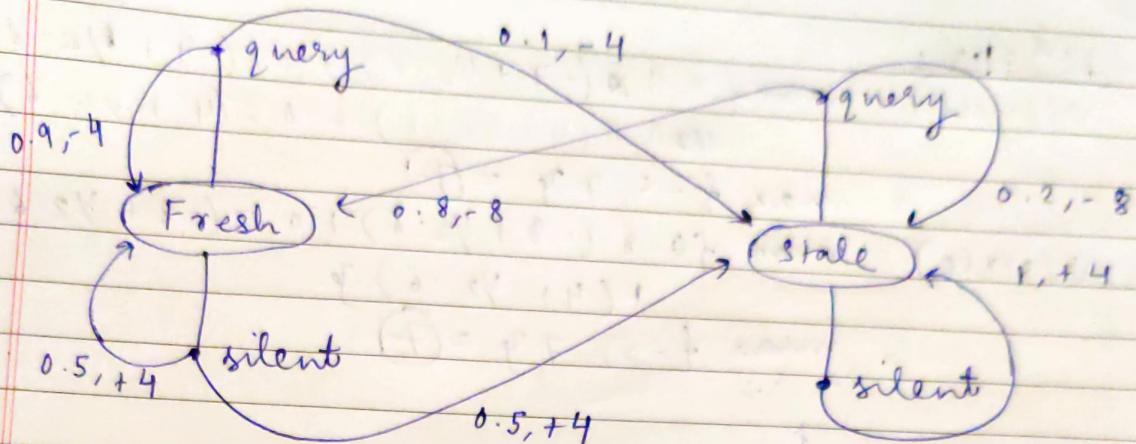
$$\begin{aligned}
 v_{\pi}(s) &\leq q_{\pi}(s, \pi'(s)) \\
 &= \mathbb{E}_{\pi}[R_{t+1} + \gamma v_{\pi}(s_{t+1}) \mid s, \pi'(s)] \\
 &= \mathbb{E}_{\pi'}[R_{t+1} + \gamma v_{\pi}(s_{t+1}) \mid s] \\
 &\leq \mathbb{E}_{\pi'}[R_{t+1} + \gamma q_{\pi}(s_{t+1}, \pi'(s_{t+1})) \mid s] \\
 &= \mathbb{E}_{\pi'}[R_{t+2} + \gamma \mathbb{E}_{\pi}[R_{t+2} + \gamma v_{\pi}(s_{t+3})] \mid s_{t+1}] \\
 &= \mathbb{E}_{\pi'}[R_{t+1} + \gamma R_{t+2} - \gamma^2 v_{\pi}(s_{t+3}) \mid s]
 \end{aligned}$$

From here

$$\begin{aligned}
 &= \mathbb{E}_{\pi'}[R_{t+1} + \gamma R_{t+2} - \gamma^5 R_{t+6} \mid s] \\
 &= \mathbb{E}_{\pi'}[G_t \mid s_t = s] \\
 &= v_{\pi'}(s)
 \end{aligned}$$

Hence, $v_{\pi}(s) \leq v_{\pi'}(s)$

16. (a)



Value Iteration

(c) $\gamma = 1/2$, $v(fresh) = 0$, $v(stale) = 0$, initially.

$$v_{k+1}(s) = \max_a \sum_{s', r} p(s'|s, a)(r + \gamma v_k(s'))$$

1st iter:

$$\begin{aligned} v(fresh) &= \max \{ 0.9(-4 + 1/2 \cdot 0) + 0.1(-4 + 1/2 \cdot 0), \\ &\quad 0.5(4 + 1/2 \cdot 0) + 0.5(4 + 1/2 \cdot 0) \} \\ &= \max \{ -4, 4 \} = 4 \end{aligned}$$

$$\begin{aligned} v(stale) &= \max \{ 0.8(-8 + 1/2 \cdot 0) + 0.2(-8 + 1/2 \cdot 0), \\ &\quad 0.5(4 + 1/2 \cdot 0) \} \\ &= \max \{ -8, 4 \} = 4 \end{aligned}$$

2nd iter:

$$\begin{aligned} v(fresh) &= \max \{ 0.9(-4 + 0.5(4)) + 0.1(-4 + 1/2 \cdot 4), \\ &\quad 0.5(4 + 1/2 \cdot 4) + 0.5(4 + 0.5(4)) \} \\ &= \max \{ -2, 6 \} = 6 \end{aligned}$$

$$\begin{aligned} v(stale) &= \max \{ 0.8(-8 + 1/2 \cdot 4) + 0.2(-8 + 1/2 \cdot 4), \\ &\quad 0.5(4 + 1/2 \cdot 4) \} \\ &= \max \{ -6, 6 \} = 6 \end{aligned}$$

3rd iter:

$$v(\text{fresh}) = \max \{ 0.9 (-4 + 1/2 \cdot 6) + 0.1 (-4 + 1/2 \cdot 6), \\ 0.5 (4 + 1/2 \cdot 6) + 0.5 (4 + 1/2 \cdot 6) \}$$

$$= \max \{ -5, 7 \} = 7$$

$$v(\text{stale}) = \max \{ 0.8 (-8 + 1/2 \cdot 6) + 0.2 (-8 + 1/2 \cdot 6), \\ 1 (4 + 1/2 \cdot 6) \}$$

$$= \max \{ -5, 7 \} = 7$$

4th iter:

$$v(\text{fresh}) = \max \{ 0.9 (-4 + 1/2 \cdot 7) + 0.1 (-4 + 1/2 \cdot 7), \\ 0.5 (4 + 1/2 \cdot 7) + 0.5 (4 + 1/2 \cdot 7) \}$$

$$= \max \{ -4.5, 7.5 \} = 7.5$$

$$v(\text{stale}) = \max \{ 0.8 (-8 + 1/2 \cdot 7) + 0.2 (-8 + 1/2 \cdot 7), \\ 1 (4 + 1/2 \cdot 7) \}$$

$$= \max \{ -4.5, 7.5 \} = 7.5$$

~~Policy Iteration~~

$$v_{KH}(s) = \sum_a \pi_K(a|s) \sum_{s', r} p(s', r | s, a) \cdot (r + \gamma v_K(s'))$$

$$\pi_{KH}(s) = \max_a \sum_{s', r} p(s', r | s, a) \cdot (r + \gamma v_K(s'))$$

After we do this, we will go onto policy improved