Priyash ____
2018MEB 1022

__Q1.__

$v \geq u$.    P.t.    $Bv \geq Bu$.

$B$ is the optimality bellman operator.

$$\therefore \quad B(v(s)) = \max_{a} \left[ \sum_{s', r} P(s', r|s, a) \left[ r + \gamma v(s') \right] \right].$$

value inside the max operator

$$\sum_{s', r} P(s', r|s, a) \left[ r + \gamma v(s') \right]$$

$$\geq \sum_{s, r} P(s', r|s, a) \left[ r + \gamma u(s') \right]$$

$$\therefore \quad v(s') \geq u(s').$$

$\therefore$ For any action the value of LHS (with $v$) $\geq$ RHS (with $u$). But this will guarantees that the best action for $v \geq$ best action for $u$.

$\therefore$ if the value of RHS with the best action '$a$' is determined, it is certain the LHS will have a value $\geq$ at that point. (or maybe some where else too then that will be the max.)

$$\therefore \quad \max_{s' r} \sum P(s', r|s, a) \left[ r + \gamma v(s') \right]$$

$$\geq$$

$$\max_{s' r} \sum P(s', r|s, a) \left[ r + \gamma u(s') \right]$$

$$\Rightarrow \quad Bv \geq Bu$$

**Q2.**

$[v_0, v^1, v^2 \ldots v^n \ldots v*] \equiv$ iterations $\theta$ value itr algorithm.

$$\therefore \|B v^{n-1} - B v*\| \leq \|v^{n-1} - v*\| \, r$$

$$\Rightarrow \|v^n - v*\| \leq r \, \|v^{n-1} - v*\|$$

$$\because \|v^{n-1} - v*\| \leq r \|v^{n-2} - v*\| \text{ (or) } v^k - v* \leq r \|v^{k-1} - v*\|$$

$$\therefore \quad \boxed{\|v^n - v*\| \leq r^n \, \|v^0 - v*\|}$$

ALSO, from $\Delta$ inequality,

$$\|v^0 - v*\| - \|v^1 - v*\| \leq \|v^0 - v^1\|$$

$$\|v^0 - v*\| \leq \|v^0 - v^1\| + \|v^1 - v*\|$$

on extending,

$$\|v^0 - v*\| \leq \|v^0 - v^1\| + \|v^1 - v^2\| + \|v^2 - v^3\|$$
$$\ldots \|v^* - v*\| \to 0$$

$$\boxed{\therefore \|v^n - v*\| \leq r^n \left(\|v^0 - v^1\| + \|v^1 - v^2\| + \|v^3 - v^2\| \ldots\right)}$$

ALSO, $\|B v^k - B v^{k-1}\| \leq r \|v^k - v^{k-1}\|$

$$\|v^{k+1} - v^k\| \leq r \|v^k - v^{k-1}\|$$

Telesopically, $\|v^{k+1} - v^k\| \leq r^k \|v^1 - v^0\|$

$$\hookrightarrow \|v^2 - v^1\| \leq r \|v^1 - v^0\|$$

$$\hookrightarrow \|v^3 - v^2\| \leq r^2 \|v^1 - v^0\| \quad \circ \circ \circ$$

$$\|v^n - v*\| \leq r^n \left(\|v^0 - v^1\| + r \|v^0 - v^1\| + r^2 \|v^0 - v^1\| + \ldots\right)$$

$$\leq r^n (v^0 - v^1)(1 + r + r^2 + \ldots)$$

$$\leq \frac{r^n}{1 - r} \|v^0 - v^1\|$$

## Q3.

### 1. T.C. of value iteration

$\Rightarrow$ $O(|S|^2 |A|)$ for each iteration, but there will be many iterations.

[ linear convergence ]

$\downarrow$

For every state the no. of new steps it can reach by taking any action $\Rightarrow$ $|S| * |A|$ (maximum).

~~Past~~ There are $|S|$ such states

$TC = \therefore O(|S| * |S| * |A|) = O(|S|^2 |A|)$

$|S| \equiv$ no of states
$|A| \equiv$ No of actions possible at any state.

### 2. Policy Iteration

$\rightarrow$ Lower # of iterations.

$\Rightarrow$ $O(|S|^3 + |S|^2 |A|)$ $\nearrow$

The $|S|^2 |A|$ comes from the value iteration part, but the policy evaluation part adds to the overall TC.

Normal policy evaluation $\Rightarrow$ $O(|S|^3)$

$\therefore$ overall TC of policy itr $\Rightarrow$ $O(|S|^3) + O(S^2 A)$

$= O(|S|^3 + |S|^2 |A|)$

### 3. Modified Policy iteration $\rightarrow$ Lower # of iterations.

$\Rightarrow$ $O(|S|^2 K + |S|^2 |A|)$

$\rightarrow$ same as policy itr except the evaluation runs for $K$ steps only. $\Rightarrow$ this step $= O(|S|K)$

There are $|S|$ such states $\Rightarrow$ $O(|S|^2 K)$

overall $\Rightarrow$ $O(|S|^2 (K + |A|))$

Q4) $q_\pi(S,a) > V_\pi(S)$

$V_\pi(S) = \sum\limits_{a'} \pi(a'|S)\, q_\pi(S,a')$ , $a' \in A(S)$.

∴ The value func. is the expected value of $q_\pi(S,a')$,
or a weighed avg of it, the weights are the distribution of
$\pi(a'|S)$ ⌡

depends on the
policy.

It is given that for some
$S \in S$ & $a \in A(S)$,

$q_\pi(S,a) > V_\pi(S)$.

⇒ This means that there exists a state$^S$, where an action$^a$
can be made that results in a higher expected reward
then the one given by the policy $\pi$.

⇒ This means there exists another Policy $\pi^{new}$
where the action taken will be more exploitory
that $\pi^{new}(a|S) > \pi(a|S)$

& will increase $V_{\pi new}(S)$.

∴ $V_\pi \neq$ an optimal policy for sure.

**Q5.** $0 < \gamma < 1$ $\&$ $\gamma \neq 1$ $\therefore$ (Horizon is inf)

Now, $\quad V_\pi(s) = E_\pi\left[G_t \mid S_t = s\right]$ $\quad , \forall s \in S$.

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

$\therefore$ if $v_\pi^{new}$, $R_{k+1}^{new} = C + R_{k+1}$

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

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

$$= V_\pi(s) + \underbrace{E_\pi\left[\sum_{k=0}^{\infty} \gamma^k C\right]}_{\text{constant.}}$$

**General.**
$$\boxed{v_\pi^{new}(s) = V_\pi(s) + \sum_{k=0}^{n} \gamma^k C}$$

for this question,
$0 < \gamma < 1$ $\&$ $n \to \infty$

$\therefore$

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

$$= V_\pi(s) + \left[C + C\gamma + C\gamma^2 + \cdots \infty\right]$$

$$= V_\pi(s) + \frac{C}{1-\gamma}$$

## Q6

**a.)** ~~shortest~~ shortest path → least amount of steps.

∴ if $R_s = -1$, when the agent will try to maximize returns, it will find the shortest path.

$$V_*(S) = \max_a \sum_{s', r} P(s', r | s, a)\left[ r + V_*(s')\right]$$

∵ Actions are deterministic, $P'(s', r | s, a) = 1$.

Given → $R_G = +5$, $\gamma = 1$.

$R_D = -5$

$R_{12} = +5$   $R_5 = -5$

$$V_*(S) = \max_a \sum_{s', r}\left[ r + V_*(s')\right]$$

$$\therefore V_*(8) = V_*(11) = \overset{max}{\;}(-1 + 5, 0, 0) = \overset{max}{\;}(4, 0, 0) = 4.$$

For ~~$V_*(7)$~~ $V_*(7)$, it would be $\overset{max}{\;}(-1 + 4, 0, 0) = 3$.

∴ It can be observed that the value of any state ⟹

$$5 - (\text{Its distance from } G = 12).$$

| | | | |
|---|---|---|---|
| $V_1 = 0$ | $V_6 = 2$ | $V_{10} = 3$ | $V_{15} = -1$ |
| $V_2 = 1$ | $V_7 = 3$ | $V_{11} = 4$ | $V_{16} = -2$ |
| $V_3 = 2$ | $V_8 = 4$ | $V_{13} = 1$ | $V_5 = 0 - 5$ |
| $V_4 = 3$ | $V_9 = 2$ | $V_{14} = 0$ | $V_{12} = 0 + 5$ |

} for optimum value function.

---

**b.)** From Q5. except for states $= \infty$, it is the same question.   [the general sol\ⁿ]

$\gamma = 1$.

$c = +2$.

$$V_\pi^{new} = V_\pi + \sum_{k=0}^{n} c \cdot 1^k$$

$$V_\pi^{new} = V_\pi + 2 \sum_{k=0}^{k=n} (1)$$

$$V_\pi^{new} = V_\pi + 2(n+1)$$

$\left[\begin{array}{l} n \equiv \text{no. of steps} \\ \text{from current} \\ \text{state to target} \\ S = 12. \end{array}\right]$

$$\Rightarrow V_\pi^{new} = V_\pi + 2(5 - V_\pi) \qquad \therefore V_\pi = +5 - n$$

$$= V_\pi + 12 - 2V_\pi$$

$$= 12 - V_\pi$$

| | | | |
|---|---|---|---|
| $V_1^{new} = 10$ | $6 = 8$ | $10 = 7$ | $15 = 11$ |
| $2 = 9$ | $7 = 7$ | $11 = 6$ | $16 = 12$ |
| $3 = 8$ | $8 = 6$ | $13 = 9$ | $V_5^{new} = -3$ |
| $4 = 7$ | $9 = 8$ | $14 = 10$ | $V_{12}^{new} = 7$ |

→ [these values $+2$.]