

Ques 1. Given $u \in V$, $v \in V$ and $v \geq u$

Let there be s states denoted by $1, \dots, s$.

Now we know that

$$v_1 > u_1$$

$$v_2 > u_2$$

\vdots

$$v_s > u_s$$

where v_i represents value corresponding to value function ' v ' and state ' i '.
Similar representation for u_i .

$$\text{Now... } B_v = \max_a \sum P(s', r | s, a) [\tau + \gamma V(s')] \dots \textcircled{1}$$

$$\text{and } B_u = \max_a \sum P(s', r | s, a) [\tau + \gamma U(s')] \dots \textcircled{2}$$

Since both u and $v \in V$... probabilities and reward will be same.

Let say that $\textcircled{2}$ is maximum when we goto state s_k
i.e. $\max_a \sum P(s', r | s, a) [\tau + \gamma U(s')]$ is max^m when s' is s_k .
(w.l.o.g we can say action taken is a_k)

But...

$$\sum P(s', r | s, a) [\tau + \gamma V(s')] > \sum P(s', r | s, a) [\tau + \gamma U(s_k)]$$

i.e. when we take action a_k and we are in state s

and by taking action a_k we reach s_k then $\textcircled{2}$ is max^m.

Also we know that for each state $v_i \geq u_i$

$\Rightarrow v_{s_k} \geq u_{s_k}$. Therefore for same action taken

value of $\textcircled{1} \geq \textcircled{2}$ Hence $B_v \geq B_u$. Note that

a_k action may not give max^m for B_v but we

are sure that if for action a_k , $\textcircled{1} \geq \textcircled{2}$ then

for $B_v \geq B_u$ will also hold.

Ques 2.

Given $\{v_0, v_1, v_2, \dots, v_n, \dots\}$ successive iteration value function.

$v^* \leftarrow$ optimal value function

$$v_1 = Bv_0$$

$$v_2 = B(v_1) = B(Bv_0) = B^2(v_0)$$

\vdots

$$v^n = B^n(v_0).$$

and

$$v^* = Bv^* = B^n v^* \quad \text{since } v^* \text{ is optimal.}$$

Now...

$$\begin{aligned} \|v^n - v^*\| &= \|B^n(v_0) - B^n(v^*)\| \\ &\leq \gamma^n \|v_0 - v^*\| \dots \textcircled{1} \end{aligned}$$

$$\begin{aligned} \|v_0 - v^*\| &\leq \|v_0 - v_1\| + \|v_1 - v_2\| + \dots \\ &\leq \|v_0 - v_1\| (1 + \gamma + \gamma^2 + \dots) \\ &\leq \frac{\|v_0 - v_1\|}{1 - \gamma} \dots \textcircled{2} \end{aligned}$$

Combining $\textcircled{1}$ and $\textcircled{2}$...

$$\|v^n - v^*\| \leq \frac{\gamma^n}{1 - \gamma} \|v_0 - v_1\|.$$

(hence proved).

Ques 3.
(a)

Runs for
 $|A| |S|^2$
times

loop:

$$\Delta \leftarrow 0$$

loop for each $s \in S$: } \leftarrow Runs for $|S|$ times

$$v \leftarrow V(s)$$

$$V(s) \leftarrow \max_a \sum_{s', r} p(s', r | s, a) [r + \gamma V(s')] \}$$

Runs for
 $|A| \times |S|$ times

$$\Delta \leftarrow \max(\Delta, |v - V(s)|)$$

until $\Delta < \epsilon$

Output a deterministic policy $\pi \approx \pi^*$, such that

$$\pi(s) = \arg \max_a \sum_{s', r} p(s', r | s, a) [r + \gamma V(s')].$$

\rightarrow

$$O(|S|^2 |A|)$$

as $|A|$ action
and we can

goto $|S|$ states
(next states)

Ques 3. Given S states, A actions.

(b)

1. Initialization

2. Evaluation $\leftarrow |S|^3$ times

loop :

$\Delta \leftarrow 0$

$S_{times} \leftarrow$ loop for each $s \in S$:

$v \leftarrow V(s)$

at most S_{times} $V(s) \leftarrow \sum_{s',r} p(s',r|s,a) [r + \gamma V(s')]$

as at most S $\Delta \leftarrow \max(\Delta, |v - V(s)|)$

next steps until $\Delta < \theta$.

\rightarrow Terminates at most S_{times}

So $\rightarrow |S| \times |S| \times |S|$ times.

S_{times}
 $\times S_{times}$

3. Improvement

policy stable \leftarrow true

for each $s \in S \leftarrow S_{times}$

old action $\leftarrow \pi(s)$

$\pi(s) \leftarrow \arg \max_a \sum_{s',r} p(s',r|s,a) [r + \gamma V(s')]$

if old action $\neq \pi(s)$, policy stable \leftarrow false

if policy stable then stop!

$|A| \times |S| \times S_{times}$
for each action
check
 \downarrow
we can go to
next S at most
state

overall time complexity : $|S|^3 + |A||S|^2$ times.

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

Ques 3.

(C)

if policy evaluation step runs
for k iterations then it will take $O(k|S|^2)$.

Policy improvement will be same i.e. $O(|S|^2|A|)$.

→ Overall : $O(k|S|^2 + |S|^2|A|)$.

Ques 4. given $q_{\pi}(s, a) > v_{\pi}(s) \dots$

Policy improvement theorem states that if $v_{\pi}(s) \leq q_{\pi}(s, \pi'(s))$
then $v_{\pi}(s) \leq v_{\pi'}(s) \dots \textcircled{1}$

here if the action a is taken by policy π' i.e. $\pi'(s) = a$
then... using policy improvement theorem...

if $\pi'(s) = a \dots \textcircled{2}$

$q_{\pi}(s, a) > v_{\pi}(s) \dots \textcircled{3}$

then $v_{\pi'}(s) > v_{\pi}(s)$ using $\textcircled{1}, \textcircled{2}$ and $\textcircled{3} \dots$

hence we can say that there exists a better policy
than π for given case i.e. π is not an optimal
policy.

~~Proof of policy imp theorem...~~

$$\begin{aligned} v_{\pi}(s) &\leq q_{\pi}(s, \pi'(s)) \quad (\text{where } \pi'(s) = a) \\ &= E [R_{t+1} + \gamma v_{\pi}(s_{t+1}) | s_t = s, A_t = \pi'(s)] \\ &= E_{\pi'} [R_{t+1} + \gamma v_{\pi}(s_{t+1}) | s_t = s] \\ &< E_{\pi'} [R_{t+1} + \gamma q_{\pi}(s_{t+1}, \pi'(s_{t+1})) | s_t = s] \\ &= E_{\pi'} [R_{t+1} + \gamma R_{t+2} + \gamma^2 v_{\pi}(s_{t+2}) | s_t = s] \\ &\vdots \\ &< E_{\pi'} [R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} \dots | s_t = s] \\ &< v_{\pi'}(s) \end{aligned}$$

Solution: Earlier... if we are in state s and we goto state s_1 to s_2

Ques 5.

$$s \xrightarrow{r_1} s_1 \xrightarrow{r_2} s_2 \dots$$

expected return will be $r_1 + \gamma r_2 + \gamma^2 r_3 \dots$ for (s)

now...

$$s \xrightarrow{r_1+c} s_1 \xrightarrow{r_2+c} s_2 \dots$$

$$\text{i.e. } (r_1 + \gamma r_2 + \dots) + (c + \cancel{\gamma c} + \gamma^2 c + \dots)$$

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

$$V_{\pi}^{\text{new}}(s) = V_{\pi}(s) + \frac{c}{1-\gamma}$$

Since states are infinite and no terminal state
this is valid for all s .

$$V_{\pi}^{\text{new}} = V_{\pi} + \frac{c}{1-\gamma}$$

Ques 6.

(a) $R_s = -1.$

$$V(1) = 0$$

$$V(2) = +1$$

$$V(3) = +2$$

$$V(4) = +3$$

$$V(5) = -5$$

$$V(6) = +2$$

$$V(7) = +3$$

$$V(8) = +4$$

$$V(9) = +2$$

$$V(10) = +3$$

$$V(11) = +4$$

$$V(12) = +5$$

$$V(13) = +1.$$

$$V(15) = -1$$

$$V(16) = -2$$

$$V(14) = 0$$

(D) If we are using ^{policy} previous, then we will still get the shortest path but value function will be different.

$$\begin{aligned}V(1) &= +12 \\V(2) &= +11 \\V(3) &= +10 \\V(4) &= +9 \\V(5) &= -3\end{aligned}$$

$$\begin{aligned}V(6) &= +10 \\V(7) &= +9 \\V(8) &= +8 \\V(9) &= +10 \\V(10) &= +9 \\V(11) &= +8\end{aligned}$$

$$\begin{aligned}V(12) &= +7 \\V(13) &= +11 \\V(14) &= +12 \\V(15) &= +13 \\V(16) &= +14 \\&v\end{aligned}$$