

## ASSIGNMENT - 2

- ① a) from the description in question we get that  $V_{k+1}$  is the ~~last~~ last iteration of the algorithm, and also

$$\|V_{k+1} - V_k\|_{\infty} \leq \epsilon$$

let  $B$  be the Bellman evaluation backup

$$V_{k+1} = BV_k$$

$$\|V_k - \bar{V}\|_{\infty} \leq \|V_k - V_{k+1}\|_{\infty} + \|V_{k+1} - \bar{V}\|_{\infty}$$

(triangular inequality over norms)

$$= \|V_k - V_{k+1}\|_{\infty} + \|BV_k - B\bar{V}\|_{\infty}$$

(Bellman evaluation backup)

$$\|V_k - \bar{V}\|_{\infty} \leq \epsilon + \gamma \|V_k - \bar{V}\|_{\infty}$$

$$\|V_k - \bar{V}\|_{\infty} \leq \frac{\epsilon}{1-\gamma}$$

$$\|V_{k+1} - \bar{V}\|_{\infty} = \|BV_k - B\bar{V}\|_{\infty} \leq \gamma \|V_k - \bar{V}\|_{\infty}$$

$$\|V_{k+1} - \bar{V}\|_{\infty} \leq \frac{\gamma \epsilon}{1-\gamma}$$



① b)

from the previous question we get that

$$\|V_{k+1} - V^\pi\|_\infty \leq \gamma \|V_k - V^\pi\|_\infty \quad \text{--- (1)}$$

similarly

$$\|V_k - V^\pi\|_\infty \leq \gamma \|V_{k-1} - V^\pi\|_\infty$$

⋮

$$\|V_2 - V^\pi\|_\infty \leq \gamma \|V_1 - V^\pi\|_\infty$$

$$\|V_{k+1} - V^\pi\|_\infty \stackrel{\text{using this in eq (1)}}{\leq} \gamma^2 \|V_{k-1} - V^\pi\|_\infty$$

$$\|V_{k+1} - V^\pi\|_\infty \leq \gamma^k \|V_1 - V^\pi\|_\infty \quad \text{similarly replacing each equation on RHS}$$

$$\|V_{k+1} - V^\pi\|_\infty \leq \gamma^k \|V_1 - V^\pi\|_\infty$$

② given that action set  $A$  is finite, we can say that there exists  $a_1, a_2 \in A$  such that for a fixed state  $s \in S$

$$L(u(s)) = \left[ R(s, a_1) + \gamma \sum_{s'} P(s'/s, a_1) u(s') \right]$$

--- (1)

and also

$$L(v(s)) = \left[ R(s, a_2) + \gamma \sum_{s'} P(s'/s, a_2) v(s') \right]$$

from the definition of optimality operator  $a_2 \in A$  gives the maximum value for  $L$ , we can say that

$$L(v(s)) \geq \left[ R(s, a_1) + \gamma \sum_{s'} P(s'/s, a_1) v(s') \right] \quad \text{--- (2)}$$

subtract (2) from (1)

$$L(u(s)) - L(v(s)) \leq \left[ \gamma \sum_{s'} P(s'/s, a_1) (u(s') - v(s')) \right]$$

since  $u(s) \leq v(s)$  ;

$$L(u(s)) - L(v(s)) \leq \left[ \gamma \sum_{s'} P(s'/s, a_1) (v(s') - v(s')) \right]$$

$$= 0$$

$$L(u(s)) - L(v(s)) \leq 0$$

$$L(u(s)) \leq L(v(s))$$

$$L(u) \leq L(v)$$

hence proved



②②

for all  $u, v \in V$ , we can write

$$\|P \circ Q(v) - P \circ Q(u)\| = \|P(Q(v)) - P(Q(u))\|$$

$$\leq r_p \|Q(v) - Q(u)\|$$

$$= r_p r_q \|v - u\|$$

$$\|P \circ Q(v) - P \circ Q(u)\| \leq r_p r_q \|v - u\|$$

$r_p$  and  $r_q$  both belong to  $[0, 1]$   
hence their product also belongs to  $[0, 1]$   
which makes  $P \circ Q$  contraction to be  
contraction on the same normed  
vector space.

$$\|Q \circ P(v) - Q \circ P(u)\| = \|Q(P(v)) - Q(P(u))\|$$

$$\leq r_q \|P(v) - P(u)\|$$

$$\|Q \circ P(v) - Q \circ P(u)\| \leq r_p r_q \|v - u\|$$

Similarly,  $Q \circ P$  contraction is also contraction  
on the same normed vector space.

②③

from the ②② we can see that  
the contraction coefficient for  $P \circ Q$  and  
 $Q \circ P$  is  $r_p r_q$ .

②④

For a unique solution to exist, the  
operator  $B$  as  $F \circ I$  must be contraction.  
which means that both  $F$  and  $I$   
must be contraction under the max-norm.



③ (a) Trajectories starting from  $s$  looks like below

$$\{ \underbrace{s, s, s, \dots, s}_K, A \}$$

K occurrences  
of  $s$

~~③ (b) First time estimate of  $v(s)$  as the reward for state is 1~~

③ (b) Depending on the reward for state  $s$  which is 1, from the state  $s$  the typical trajectory has  $k$  number of occurrences of  $s$  then the first visit estimate of  $v(s)$  is  $k$ .

③ (c) Given one trajectory, we can construct  $k-1$  sub-trajectories (no. of times state  $s$  is visited) where  $t$ th trajectory returns  $k-t$ . ( $k$  is the no. of occurrences of state  $s$ )

$$\hat{v}(s) \text{ (estimated } v) = \frac{1}{k} \sum_{t=0}^{k-1} (k-t)$$

$$= \frac{1}{k} \sum_{t'=1}^k t'$$

$$= \frac{1}{k} \left( \frac{k(k+1)}{2} \right)$$

$$\hat{v}(s) = \frac{k+1}{2}$$

③ (d)

$$v(s) = 1 + (1-\beta) v(s) + \beta \cdot p$$

$$v(s) = \frac{1 + \beta v(s)}{1 - (1-\beta)}$$

$$v(s) = \frac{1}{\beta}$$

Page \_\_\_\_\_

$$\mathbb{E}[\hat{v}] = \mathbb{E}[\underbrace{K}_{\substack{\text{k occurrences} \\ \text{of state s}}}] = \frac{1}{p} = v^{\pi}(s)$$

So the true value of  $v(s)$  is 1 as its average length of the trajectory

③② The value function  $\hat{v}(s)$  estimation for every visit is given as  $\frac{R+1}{2}$

The corresponding expectation can be written as

$$\begin{aligned} \mathbb{E}\left[\frac{R+1}{2}\right] &= \frac{1}{2} \mathbb{E}[K] + \frac{1}{2} \\ &= \frac{1}{2} \left[ \frac{1}{p} + 1 \right] \\ &= \frac{1+p}{2p} \neq v^{\pi}(s) \end{aligned}$$

This means that the every visit MC estimate is biased estimator.

③③ - First visit MC converges just rely on the law of large numbers whereas every visit MC convergence may need to more care because the samples for mean calculation are all independent.



$$(4) (a) \quad \delta_t = r_{t+1} + \gamma v^\pi(s_{t+1}) - v^\pi(s_t)$$

$$\mathbb{E}[\delta_t | s_t = s] = \mathbb{E}[r_{t+1} + \gamma v^\pi(s_{t+1}) - v^\pi(s_t) | s_t = s]$$

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

$$= v^\pi(s) - v^\pi(s) \\ \mathbb{E}[\delta_t | s_t = s] = 0$$

$$(4) (b) \quad \mathbb{E}[\delta_t | s_t = s, A_t = a]$$

$$= \mathbb{E}[r_{t+1} + \gamma v^\pi(s_{t+1}) - v^\pi(s_t) | s_t = s, A_t = a]$$

$$= \mathbb{E}[r_{t+1} + \gamma v^\pi(s_{t+1}) | s_t = s, A_t = a] - v^\pi(s)$$

$$= Q^\pi(s, a) - v^\pi(s)$$

$$\mathbb{E}[\delta_t | s_t = s, A_t = a] = Q^\pi(s, a) - v^\pi(s)$$

(4) (c)  $\eta(\lambda)$  denote the time by which the weighting sequence would have fallen half of its initial value. The half life occurs when weighting drops by half. Therefore

$$\lambda^{\eta(\lambda)} = \frac{1}{2}$$

$$\eta(\lambda) = \frac{\ln(\frac{1}{2})}{\ln(\lambda)}$$

for 3 step return, the time taken is 2

$$\eta(\lambda) = 2 \Rightarrow 2 = \frac{\ln(\frac{1}{2})}{\ln(\lambda)}$$

$$\Rightarrow \lambda^2 = \frac{1}{2}$$

$$\lambda = \frac{1}{\sqrt{2}}$$

⑤ The harmonic series' generalization is the  $p$ -series defined as  $\sum_{i=1}^{\infty} \frac{1}{i^p}$

$$\sum_{i=1}^{\infty} \frac{1}{i^p} \quad \text{for any postive real number } p$$

The  $p$ -series converges for all  $p > 1$  and diverges for all  $p \leq 1$ .

$$(i) \quad \alpha_t = \frac{1}{t} \Rightarrow \sum_{t=1}^{\infty} \alpha_t = \infty$$

$$\sum_{t=1}^{\infty} \alpha_t^2 < \infty$$

Algo doesnot converge with given  $\alpha_t$ .

$$(ii) \quad \alpha_t = \frac{1}{t^2} \Rightarrow \sum_{t=1}^{\infty} \alpha_t = \infty$$

$$\sum_{t=1}^{\infty} \alpha_t^2 < \infty$$



Algo converges with given value of  $\alpha_t$ .

$$(iii) \quad \alpha_t = \frac{1}{t^{2/3}} \Rightarrow \sum_{t=1}^{\infty} \alpha_t = \infty$$

$$\sum_{t=1}^{\infty} \alpha_t^2 < \infty$$

Algo converges with given value of  $\alpha_t$ .

$$(iv) \quad \alpha_t = \frac{1}{t^{1/2}} \Rightarrow \sum_{t=1}^{\infty} \alpha_t = \infty$$

$$\sum_{t=1}^{\infty} \alpha_t^2 = \infty$$

do not  
Algo converges with given value of  $\alpha_t$ .

generalizing above result we can write

$$\sum_{t=1}^{\infty} \frac{1}{t^p} = \infty \quad p \leq 1$$

$$\sum_{t=1}^{\infty} \frac{1}{t^{2p}} < \infty \quad \begin{matrix} 2p > 1 \\ p > \frac{1}{2} \end{matrix}$$

The Robbins - Monro condition is true when  $p \in (\frac{1}{2}, 1)$ .