

Reinforcement - Learning

Assignment - 2

Question - 1)

a) we terminate the algorithm when

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

Lets say we didn't stop it at this condition, instead we let the value function converge.

The iteration scheme is

$$V^{(k+1)}(s) = \sum_a \pi(a|s) \sum_{s'} P_{ss'}^a (R_{ss'}^a + \gamma V^{(k)}(s)) \quad \text{--- (1)}$$

$$V^{(k+2)}(s) = \sum_a \pi(a|s) \sum_{s'} P_{ss'}^a (R_{ss'}^a + \gamma V^{(k+1)}(s)) \quad \text{--- (2)}$$

$$\text{(2)} - \text{(1)}$$

$$V^{(k+2)}(s) - V^{(k+1)}(s) = \sum_a \pi(a|s) \sum_{s'} P_{ss'}^a \gamma (V^{(k+1)}(s) - V^{(k)}(s)) \quad \text{--- (3)}$$

we know that $\|V^{(k+1)} - V^{(k)}\|_{\infty} \leq \epsilon$

$$\Rightarrow |V^{(k+1)}(s) - V^{(k)}(s)| \leq \epsilon \quad \forall s \in S$$

$$\Rightarrow -\epsilon \leq V^{(k+1)}(s) - V^{(k)}(s) \leq \epsilon \quad \text{--- (4)}$$

Using (3) and (4)

$$|V^{(k+2)}(s) - V^{(k+1)}(s)| \leq \sum_a \pi(a|s) \gamma \epsilon \quad \forall s$$

$$\boxed{-\gamma \epsilon \leq V^{(k+2)}(s) - V^{(k+1)}(s) \leq \gamma \epsilon}$$

similarly we can say

$$-\gamma^2 \epsilon \leq v^{(k+3)} - v^{(k+2)}(s) \leq \gamma^2 \epsilon$$

Now

$$-\gamma \epsilon \leq v^{(k+2)}(s) - v^{(k+1)}(s) \leq \gamma \epsilon$$

$$-\gamma^2 \epsilon \leq v^{(k+3)}(s) - v^{(k+2)}(s) \leq \gamma^2 \epsilon$$

⋮

Adding these equations

$$-(\gamma \epsilon + \gamma^2 \epsilon + \dots) \leq v^\pi(s) - v^{(k+1)}(s) \leq \gamma \epsilon + \gamma^2 \epsilon + \dots$$

$$|v^\pi(s) - v^{(k+1)}(s)| \leq \gamma \epsilon + \gamma^2 \epsilon + \dots \leq \frac{\gamma \epsilon}{1 - \gamma}$$

$$\Rightarrow \boxed{\|v^\pi - v^{(k+1)}\|_\infty \leq \frac{\gamma \epsilon}{1 - \gamma}}$$

Hence Proved

b) we ~~the~~ know that L

(Bellman optimality operator)
is a contraction map.

Since v^π is the optimal value,

$$L(v^\pi) = v^\pi \iff v^\pi \text{ is fixed for}$$

from the properties of contraction
map we know that

$$\|L(v^{(k)}) - L(v^\pi)\|_\infty \leq \gamma \|v^{(k)} - v^\pi\|_\infty \quad \textcircled{1}$$

from our algorithm

$$L(v^{(k)}) = v^{(k+1)}$$

Substituting in $\textcircled{1}$ & using $L(v^\pi) = v^\pi$

$$\|v^{(k+1)} - v^\pi\|_\infty \leq \gamma \|v^{(k)} - v^\pi\|_\infty$$

$$\|v^{(k)} - v^\pi\|_\infty \leq \gamma \|v^{(k-1)} - v^\pi\|_\infty$$

\vdots

$$\|v^{(2)} - v^\pi\|_\infty \leq \gamma \|v^{(1)} - v^\pi\|_\infty$$

$$\Rightarrow \|v^{(k+1)} - v^\pi\|_\infty \leq \gamma^k \|v^{(1)} - v^\pi\|_\infty$$

$$c) \quad L(v) = \max_{a \in A} [R^a + \gamma P^a v]$$

P^a is probability ~~vector~~ ^{matrix} which has all positive entries.

~~we know that $P^a x \geq 0$ if $x \geq 0$~~

Claim: $P^a x \geq 0$ if $x \geq 0$

Proof all the entries of P are positive, and also all the entries of x are positive, we can take inner product in rowwise manner all the terms would again be non negative.

$$\Rightarrow P^a(v - u) \geq 0$$

$$\Rightarrow P^a v \geq P^a u$$

adding R^a on both sides

$$P^a v + R^a \geq P^a u + R^a$$

This holds for all a

$$\therefore \max_a (P^a v + R^a) \geq \max_a (P^a u + R^a)$$

$$\Rightarrow L(v) \geq L(u)$$

Hence Proved

Problem - 3)

a) A generic form of the trajectory is

$$\underbrace{S \dots A}$$

S can occur ≥ 1 times

b) First visit MC makes sense only when we are given some experience

However, we can find the value vector by calculating expectation of G_t

For a trajectory ~~with m~~ of size m , the state S would have been visited $m-1$ times

$$G_t^{(m)} = m-1 \quad \text{since state } S \text{ is visited } m-1 \text{ times}$$

Note that the trajectory starts with S

$$E[G_t] = \sum_{m=2}^{\infty} p(1-p)^{m-2} (m-1)$$

$$= p + 2p(1-p) + 3p(1-p)^2 + \dots$$

$$V(S) = p + 2p(1-p) + 3p(1-p)^2 + \dots$$

$$(1-p)V(S) = p(1-p) + 2p(1-p)^2 + \dots$$

$$pV(S) = p + p(1-p) + p(1-p)^2 + \dots$$

$$pV(s) = \frac{p}{p}$$

$$V(s) = 1$$

for the value at A. All the trajectories end with A and that's the first visit to A in that trajectory. Since the reward for being in state A is 0

$$V(A) = 0$$

$$\therefore V = \begin{pmatrix} \frac{1}{p} \\ 0 \end{pmatrix}$$

c) Again consider a trajectory

$$G_t^{(m)} = \frac{1 + 2 + \dots + (m-1)}{m-1} = \frac{(m-1) \cancel{m} m}{2(m-1)} = \frac{m}{2}$$

Again we are considering the trajectories with state S for every visit MC since $V(A) = 0$ using the same argument is the previous subpart.

$$E[G_t] = V(s) = \sum_{m=2}^{\infty} p(1-p)^{m-2} \frac{m}{2}$$

~~$$= p(1-p)$$~~

$$V(s) = p + p(1-p)\frac{3}{2} + p(1-p)^2 2 + \dots$$

$$V(s)(1-p) = p(1-p) + p(1-p)^2 \frac{3}{2} + \dots$$

$$V(s)p = p + \frac{1}{2}(p(1-p) + p(1-p)^2 + \dots)$$

$$= p + \frac{1}{2} \cancel{p} \frac{(1-p)}{\cancel{p}}$$

$$= p + \frac{1}{2} - \frac{p}{2}$$

$$V(s) = \cancel{\frac{1+p}{2p}} \frac{p+1}{2p} = \frac{1}{2} + \frac{1}{2p}$$

d) Using the exact method

$$V = (I - \gamma P)^{-1} R$$

$$= \left(\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} - \begin{bmatrix} p & p \\ 0 & 0 \end{bmatrix} \right)^{-1} \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

$$= \begin{bmatrix} p & -p \\ 0 & 1 \end{bmatrix}^{-1} \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

$$= \frac{1}{p} \begin{bmatrix} 1 & p \\ 0 & p \end{bmatrix} \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

$$= \begin{bmatrix} \frac{1}{p} \\ 0 \end{bmatrix}$$

e) Yes, every visit MC estimate is biased.

The true value function evaluated at state S is $\frac{1}{p}$.

while every visit MC gives us $\frac{1}{2} + \frac{1}{2p}$

we know that $p < 1$

$$\frac{1}{2p} > \frac{1}{2}$$

$$\frac{1}{p} - \frac{1}{2p} > \frac{1}{2}$$

$$\Rightarrow \frac{1}{p} > \frac{1}{2} + \frac{1}{2p}$$

\therefore there is bias

f) Every visit MC will converge faster because the number of samples would be far more in case of every visit MC compared to first visit MC.

Problem - 4

a)

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

Expectation is linear

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

$$\Rightarrow V^{\pi}(s) - V^{\pi}(s)$$

$$\boxed{= 0}$$

b)

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

$$= \underbrace{E_{\pi}[r_{t+1} + \gamma V^{\pi}(s_{t+1}) | s_t = s, a_t = a]}_{Q^{\pi}(s, a)} - \underbrace{E[V^{\pi}(s_t) | s_t = s, a_t = a]}_{V^{\pi}(s)}$$

$$= Q^{\pi}(s, a) - V^{\pi}(s)$$

c)

$$G_t^{\lambda} = (1-\lambda) \sum_{n=1}^{\infty} \lambda^{n-1} G_t^{(n)}$$

The first coefficient is

$$(1-\lambda)$$

~~The coefficient~~ The subsequent coefficients are

$$(1-\lambda)\lambda, (1-\lambda)\lambda^2, \dots$$

by definition

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

$$\eta(\lambda) \log \lambda \leq -\log 2$$

$$\eta(\lambda) \leq \frac{-\log 2}{\log \lambda}$$

we can write

$$\eta(\lambda) \leq \log_{1/\lambda} 2$$

given $\eta(\lambda) = 3$

$$3 \leq \log_{1/\lambda} 2$$

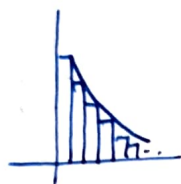
$$\lambda^{-3} \geq 2$$

$$\boxed{\lambda \geq 2^{-1/3}}$$

Problem 5

$$\alpha_t = \frac{1}{t^p}$$

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



$$\int_1^{\infty} \frac{1}{t^p} dt \leq \sum_{t=1}^{\infty} \frac{1}{t^p} \leq 1 + \int_1^{\infty} \frac{1}{t^p} dt$$

when $p > 1$

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

when $p < 1$, the summation diverges

when $p = 1$

$$\log t \Big|_1^{\infty} \leq \sum_{t=1}^{\infty} \frac{1}{t} \leq 1 + \log t \Big|_1^{\infty}$$

the summation diverges

Summary

$$\sum \frac{1}{t^p} \begin{cases} \text{converges} & p > 1 \\ \text{diverges} & p \leq 1 \end{cases}$$

$$1) \alpha_t = \frac{1}{t}$$

$$\sum \alpha_t = \infty \quad \} p=1 \text{ diverges}$$

$$\sum \alpha_t^2 < \infty \quad \} p=2 \text{ converges}$$

\therefore it would converge

$$2) \alpha_t = \frac{1}{t^2}$$

$$\sum \alpha_t = \sum \frac{1}{t^2} < \infty \quad \} p=2 \text{ converges}$$

$$\sum \alpha_t^2 = \sum \frac{1}{t^4} < \infty \quad \} p=4 \text{ converges}$$

\therefore it wouldn't converge

$$3) \alpha_t = \frac{1}{t^{2/3}}$$

$$\sum \alpha_t = \sum \frac{1}{t^{2/3}} = \infty \quad \} p=2/3 \text{ diverges}$$

$$\sum \alpha_t^2 = \sum \frac{1}{t^{4/3}} < \infty \quad \} p=4/3 \text{ converges}$$

\therefore it would converge

$$4) \quad \alpha_t = \frac{1}{t^{1/2}}$$

$$\sum \alpha_t = \sum \frac{1}{t^{1/2}} = \infty \quad \left\{ p = \frac{1}{2} \text{ diverges} \right.$$

$$\sum \alpha_t^2 = \sum \frac{1}{t} = \infty \quad \left\{ p = 1 \text{ diverges} \right.$$

\therefore it won't converge