VAISHNAVI
AI20BTECH11025

1a.   $\| V_{k+1} - V_k \|_\infty \le \epsilon$ for $\epsilon > 0$

Consider  $\| V_k - V^\pi \|_\infty$

$\| V_k - V^\pi \|_\infty \le \underbrace{\| V_k - V_{k+1} \|_\infty}_{\epsilon} + \| V_{k+1} - V^\pi \|_\infty$

$\le \epsilon + \| (R + \gamma P V_k) - (R + \gamma P V^\pi) \|_\infty$

$\le \epsilon + \gamma \| P (V_k - V^\pi) \|$

$\qquad\qquad \le \| P \| \cdot \| V_k - V^\pi \|$  $\quad$ [ $\gamma \|P\|\le$ ... ]

[∵ $\| AB \| \le \| A \| \| B \|$ ]

$\le \epsilon + \gamma \| V_k - V^\pi \| \qquad \underset{\le 1}{}$

$\Rightarrow \quad \| V_k - V^\pi \| \le \dfrac{\epsilon}{1 - \gamma}$

Consider  $\| V_{k+1} - V^\pi \|_\infty = \| (R + \gamma P V_k) - (R + \gamma P V^\pi) \|_\infty$

$\le \gamma \| P (V_k - V^\pi) \|$

$\qquad\qquad \le \underset{\le 1}{\| P \|} \cdot \| V_k - V^\pi \|$

$\le \gamma \| V_k - V^\pi \|_\infty$

$= \dfrac{\gamma \epsilon}{1 - \gamma}$  //

1b.   $\| V_{k+1} - V^\pi \|_\infty = \| (R + \gamma P V_k) - (R + \gamma P V^\pi) \|_\infty$

$\qquad\qquad\qquad \le \gamma \| V_k - V^\pi \|_\infty$

$\text{lly} \qquad \| V_k - V^\pi \| \le \gamma \| V_{k-1} - V^\pi \|_\infty$

$\qquad\qquad\qquad \vdots$

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

$\Rightarrow \quad \| V_{k+1} - V^\pi \| \le \gamma^k \| V_1 - V^\pi \|$  //

**1c.**  Bellman optimality operator

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

For value function $u$ let $a_1$ - optimal action at
$\quad\quad s$ - state

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

Value function $v$ $\quad\quad a_2$ - Optimal action at state $s$

$$L(v) = \left[ R(s, a_2) + \gamma \sum_{s'} P(s'| s, a_2) \cdot v(s) \right]$$
$$\geq \left[ R(s, a_1) + \gamma \sum_{s'} P(s'| s, a_1) v(s) \right]$$

$\Rightarrow L(u) - L(v) \leq \gamma \underset{\underset{\geq 0}{\downarrow s'}}{\sum} \underset{\geq 0}{P(s'| s, a_1)} \underset{\leq 0}{\underbrace{[u(s) - v(s)]}}$

$$\left[ \because u(s) \leq v(s) \;(\text{given}) \right.$$
$$\left. \Rightarrow u(s) - v(s) \leq 0 \right]$$

$\Rightarrow \quad L(u) - L(v) \leq 0$

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

Hence, bellman optimality operator $(L)$ satisfies the
monotonicity property.

## 2. On Contractions

a) $\nexists$ $P, Q$ are contractions on normed vector space $\langle V, \|.\|\rangle$

$\Rightarrow \|R(u) - R(v)\| \leq Q$ ~~⇒ ⅼ R(u)-R(v)∥≤ g~~

$\Rightarrow$ $\exists$ $\gamma_p, \gamma_q \in [0,1)$ s.t

$$\|P(u) - P(v)\| \leq \gamma_p \|u - v\|$$
$$\|Q(u) - Q(v)\| \leq \gamma_q \|u - v\| \qquad \forall u, v \text{ in } V$$

— Composition $P \circ Q$

$$\|P \circ Q(u) - P \circ Q(v)\| \leq \gamma_p \|Q(u) - Q(v)\|$$
$$\leq \gamma_p \cdot \gamma_q \|u - v\|$$
$$\underset{\in [0,1)}{\downarrow} \quad \underset{\in [0,1)}{\downarrow}$$
$$\Rightarrow \gamma_p \cdot \gamma_q \in [0,1)$$

$\therefore$ $P \circ Q$ is a contraction in the same normed vector space

— Composition $Q \circ P$

$$\|Q \circ P(u) - Q \circ P(v)\| \leq \gamma_q \|P(u) - P(v)\|$$
$$\leq \gamma_q \cdot \gamma_p \|u - v\|$$
$$\underset{\in [0,1)}{\downarrow} \quad \underset{\in [0,1)}{\downarrow}$$
$$\Rightarrow \gamma_q \cdot \gamma_p \in [0,1)$$

$\therefore$ $Q \circ P$ is a contraction in the same normed vector space.

b) From above, we have

Composition $P \circ Q$ : $\quad \|P \circ Q(u) - P \circ Q(v)\| \leq \gamma_p \cdot \gamma_q \|u - v\|$

$\Rightarrow \gamma_p \cdot \gamma_q$ is the suitable Lipschitz coeff.

Composition $Q \circ P$ : $\quad \|Q \circ P(u) - Q \circ P(v)\| \leq \gamma_p \cdot \gamma_q \|u - v\|$

$\Rightarrow \gamma_p \gamma_q$ is the suitable Lipschitz coeff.

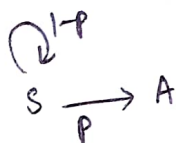2c    For the operator to converge to a unique solution, the operator FoL must be a contraction.

From (a), we have that a composition is a contraction if both the functions F & L are contractions.

$\Rightarrow$ FoL converges to unique sol when F, L are contractions under max-norm.

20a.

Q3.

$$S \xrightarrow[P]{Q^{1-p}} A$$

3a. Typical trajectory starting at $S$ : $SS \cdots A$

b. For a trajectory of length $l+1$ ie $l-'s's$

$V(S)$ first visit $MC = l$
estimate

c. $V(S)$ every visit $MC = \dfrac{l+l-1+l-2+\cdots 1}{l}$
estimate

$$= \dfrac{l(l+1)}{2l} = \dfrac{l+1}{2}$$

d. $V(A) = 0$

$V(S) = 1 + p \cdot V(A) + (1-p) V(S)$

$\Rightarrow V(S) [1 - 1 + p] = 1$

$V(S) = \dfrac{1}{P}$

e. Every visit $MC$ is biased as not all returns are not iid

Proof :

expected length of episode $= p + (1-p)p(2) + (1-p)^2 p(3) + \cdots$

$= p[1 + (1-p)2 + (1-p)^2 3 + \cdots]$

$= p\left[\dfrac{1}{p} + \dfrac{1(1-p)}{p^2}\right] = \dfrac{1}{p}$

$V(S)$ for every visit $= \dfrac{\left(\dfrac{1}{p} + 1\right)}{2}$ which is $\dfrac{1}{2}$ times that of true value.

$\Rightarrow$ Every vint $MC$ is biased.

**3f.** First visit MC: All the returns used in the calculation of value of a state are from diff episodes sampled randomly. i.e i.i.ds. By the law of large numbers it converges to the true value as the num of episodes increases.

Every visit MC: Again assuming large num of episodes and exploring different starts garantees the converges of the algorithm. Though the returns are not all iids, the bias decreases consistently with increasing num of episodes ( it asymptotically goes to zero)

**Q4.** Temporal Difference Methods

MDP - M
Policy - $\pi$

One step TD error: $\delta_t = \mathscr{R}_{t+1} + \gamma V^\pi(S_{t+1}) - V^\pi(S_t)$

a) $E_\pi(\delta_t \mid S_t = S) = E\left(\mathscr{R}_{t+1} + \gamma V^\pi(S_{t+1}) - V^\pi(S_t) \mid S_t = S\right)$

From linearity of expectation

$= E_\pi\left(\mathscr{R}_{t+1} + \gamma V^\pi(S_{t+1}) \mid S_t = S\right) - \underbrace{E\left[V^\pi(S_t) \mid S_t = S\right]}_{\displaystyle \| \atop V^\pi(S)}$

$\left[\because \text{ we are using true state value function } V^\pi\right]$

$= \sum_{a} \pi(a|s) \sum_{s'} \underbrace{P^a_{ss'}\left[R^a_{ss'} + \gamma V^\pi(s')\right]}_{\displaystyle = V^\pi(S)} - V^\pi(s)$

$\quad\quad (\text{from def of } V^\pi(s))$

$= 0$

b) $E_\pi(\delta_t \mid S_t = S, a_t = a) = E_\pi\left(\mathscr{R}_{t+1} + \gamma V^\pi(S_{t+1}) - V^\pi(S_t) \mid S_t = S, a_t = a\right)$

$= E_\pi\left[\mathscr{R}_{t+1} + \gamma V^\pi(S_{t+1}) \mid S_t = S, a_t = a\right] - \underbrace{E\left[V^\pi(S_t) \mid S_t = S, a_t = a\right]}_{\displaystyle \| \atop V^\pi(S)}$

$= \underbrace{\sum_{s'} P^a_{ss'}\left[R^a_{ss'} + \gamma V^\pi(s')\right]}_{\displaystyle Q(S,a)} - V^\pi(s)$

$\quad\quad\quad\quad [\text{from the def of } Q(S,a)]$

$= Q^\pi(S,a) - V^\pi(S)$

c) TD($\lambda$) algorithm, $\lambda$ return target

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

where $G_t^{(n)} = \mathscr{R}_{t+1} + \gamma \mathscr{R}_{t+2} + \gamma^2 \mathscr{R}_{t+3} + \dots \gamma^{n-1}\mathscr{R}_{t+n} + \gamma^n V(S_{t+n})$

Let $n(\lambda)$ denote time by which weighing seq reduces by half.

$\Rightarrow \dfrac{1}{2} \leq \dfrac{(1-\lambda)\,\lambda^{n(\lambda)-1}}{(1-\lambda)\,\lambda^{1-1}}$

$- \ln 2 \leq [n(\lambda)-1]\ln\lambda$

$\boxed{\dfrac{1-\ln 2}{\ln \lambda} \leq n(\lambda)}$

Value of $\lambda$ for which wts drop to half after 3 steps

$$\text{i.e.} \quad \eta(\lambda) = 3$$

$$1 - \frac{\ln 2}{\ln \lambda} = 3$$

$$-2 = \frac{\ln 2}{\ln \lambda}$$

$$\ln \lambda = \ln 2^{-\frac{1}{2}}$$

$$\lambda = \frac{1}{\sqrt{2}} \, \text{//}$$

Q5. Consider the p-series $\sum\limits_{n=1}^{\infty} \frac{1}{n^p}$

We know that if $\quad p \leq 1 \quad \rightarrow$ divergent

$\qquad\qquad\qquad p > 1 \quad \rightarrow$ convergent

Proof: Let us look at the convergence of the corres. integral

$$(p > 0) \quad \int\limits_{1}^{\infty} \frac{1}{n^p} \, dx = \lim\limits_{m \to \infty} \int\limits_{1}^{m} \frac{1}{n^p}$$

$p = 1$

$$\lim\limits_{m \to \infty} \left[ \ln n \right]_{1}^{M} = \ln \infty - 0$$

$$= \infty$$

$\downarrow$

Diverges

$p \neq 1$

$$\lim\limits_{m \to \infty} \left[ \frac{x^{-p+1}}{-p+1} \right]_{1}^{m}$$

$$\frac{1}{1-p} \lim\limits_{m \to \infty} m^{1-p} - \frac{1}{1-p}$$

$1 - p > 0$
$1 > p$

Diverges

$1 - p < 0$
$1 < p$

Converges

(1) $\quad a_t = \frac{1}{t}$

$$\sum\limits_{t=0}^{\infty} a_t = \sum\limits_{t=0}^{\infty} \frac{1}{t} \quad \rightarrow \quad \text{same} \quad p = 1 \rightarrow \text{diverges} = \infty$$

$$\sum a_t^2 = \sum\limits_{t=0}^{\infty} \frac{1}{t^2} \qquad p = 2 \rightarrow \text{converges} < \infty$$

$\Rightarrow$ It obeys Robbins-Monroe condition thus converges

(2)  $\alpha_t = \dfrac{1}{t^2}$

$\sum \alpha_t = \sum \dfrac{1}{t^2}$   $p > 1$  Converges $< \infty$

$\sum \alpha_t^2 = \sum \dfrac{1}{t^4}$   $p = 4 > 1$  Converges $< \infty$

$\Rightarrow$ Doesn't obey Robbins-Monroe condition thus doesn't converge.

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

$\sum \alpha_t = \sum \dfrac{1}{t^{2/3}}$   $p = 2/3 < 1$  Diverges $= \infty$

$\sum \alpha_t^2 = \sum \dfrac{1}{t^{4/3}}$   $p = 4/3 > 1$  Converges $< \infty$

$\Rightarrow$ Obeys Robbins-Monroe condition thus converges.

(4)  $\alpha_t = \dfrac{1}{t^{1/2}}$

$\sum \alpha_t = \sum \dfrac{1}{t^{1/2}}$   $p = 1/2 < 1$  Diverges $= \infty$

$\sum \alpha_t^2 = \sum \dfrac{1}{t}$   $p = 1$  Diverges $= \infty$

$\Rightarrow$ Doesn't obey Robbins-Monroe condition thus doesn't converge.

For learning rate  $\alpha_t = \dfrac{1}{t^p}$

$\sum \alpha_t = \sum \dfrac{1}{t^p} = \infty$   $\Rightarrow$  $p \le 1$

$\sum \alpha_t^2 = \sum \dfrac{1}{t^{2p}} < \infty$   $\Rightarrow$  $2p > 1$

$p > 1/2$

$\Rightarrow$  $p \in (1/2, 1]$  for $\alpha_t$ to converge to $V(s)$.