# INDIAN INSTITUTE OF TECHNOLOGY MADRAS

B

Roll No. | M | M | 1 | 6 | B | 0 | 2 | 3 |

Name : H. Vishal

Total No. of Pages

Quiz I ☐   Quiz-II/ Mid-Sem ✓   End-Semester ☐   Make-up ☐   Date : 29th Feb, 2020

Semester & Degree : 1000 DS 8th Sem   Course No. CS 6700   Part :

| Question No. | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|---|---|---|---|---|---|---|---|---|---|---|
| Marks | 0 | 0 | 6 | Ø | 2 | 0 | 0 | 1 | 1·5 | 0 |

| 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | Total |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | | | | 12·5 |

Answer on both sides of the paper including the space below

③ (a)

$$V(s) = \theta^T \phi(s) \quad \text{(Linear } f^n \text{ approximator)}$$

where $\phi(s) = \begin{bmatrix} 1 \\ -1 \\ -1 \end{bmatrix}$, $\phi(s_2) = \begin{bmatrix} -1 \\ -1 \\ 1 \end{bmatrix}$, $\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \end{bmatrix}$

$$\Rightarrow V(s) = \begin{bmatrix} V(s_1) \\ V(s_2) \end{bmatrix} = \begin{bmatrix} \theta_0 - \theta_1 - \theta_2 \\ -\theta_0 - \theta_1 + \theta_2 \end{bmatrix}$$

Consider $\theta_0 - \theta_2 = a$

$$\Rightarrow V(s) = \begin{bmatrix} a - \theta_1 \\ -(a + \theta_1) \end{bmatrix} \checkmark$$

The terms $a - \theta_1$, $-(a + \theta_1)$ are linearly independent as $\alpha(a - \theta_1) + \beta(-a - \theta_1) = 0$

$$\Rightarrow \alpha = 0, \quad \beta = 0$$

$\Rightarrow$ rank $(V) = 2$, $\dim (V) = 2$

$\Rightarrow$ Any point in $2D$ space ($\mathbb{R}^2$) can be represented and the corresponding value function can be learnt.

This is similar to a look-up table, where instead of hot-encodings, we use a different formulation.

b)     TD $(0)$ :

$$\theta_{t+1} \leftarrow \theta_t + \alpha \left( R_{t+1} + \gamma V_t(S_{t+1}) - V_t(S_t) \right) \nabla_{\theta_t} V_t(S_t)$$

$\Rightarrow$
$$\theta_{t+1} \leftarrow \theta_t + \alpha \left( -5 + \gamma V_t(S_1) - V_t(S_2) \right) \nabla_{\theta_t} V_t(S_2)$$

$$V(S_2) = \theta \phi(S_2)$$

$$\nabla_{\theta_t} V_t(S_t) = \phi(S_2) = \begin{bmatrix} -1 \\ 1 \\ 1 \end{bmatrix}$$

$\Rightarrow$
$$\theta_{t+1}^{[1]} \leftarrow \theta_t^{[1]} - \alpha \left( -5 + \gamma V_t(S_1) - V_t(S_2) \right)$$

$$\theta_{t+2}^{[2]} \leftarrow \theta_t^{[2]} - \alpha \left( -5 + \gamma V_t(S_1) - V_t(S_2) \right) \checkmark$$

$$\theta_{t+3}^{[3]} \leftarrow \theta_t^{[3]} + \alpha \left( -5 + \gamma V_t(S_1) - V_t(S_2) \right)$$

where     $\theta_t^{[i]}$     denotes the $i^{th}$ element of $\theta_t$

(2)

In case of policy improvement,

$$\pi_{n+1} \in \text{arg max}_{\pi} \left\{ r_\pi + \gamma P_\pi v^{\pi_n} \right\}$$

This is equivalent to saying

$$\pi^*(s) = \text{arg max}_a q^{\pi}(s,a)$$

$$\Rightarrow \pi^*(s) = \max_\pi \sum_{s'} p(s'|s,a) \left[ E[r|s,a,s'] + \gamma v^{\pi}(s') \right]$$

The max can be taken inside

Also, if there exists a state $s$ s.t. $q^\pi(s,a) > v^\pi(s)$

$\Rightarrow$ choosing that state would $\uparrow v^\pi$ contradicting the fact that $\pi$ is optimal

a) ~~False~~

(5)

$$G_t = R_{t+1} + R_{t+2} + \cdots R_{t+k}$$

$$v^\pi(s) = E_\pi \left[ G_t \mid S_t = s \right]$$

$$= E_\pi \left[ \left[ R_{t+1} + R_{t+2} + \cdots R_{t+k} \right] \mid S_t = s \right]$$

$$= E_\pi \left[ R_{t+1} + v^\pi(s') \mid S_t = s \right]$$
$$- R_{t+k+1}$$

$$= E_\pi \left[ (R_{t+1} - R_{t+k+1}) + v^\pi(s') \mid S_t = s \right]$$

$$\Rightarrow \quad v^*(s) = \max_a \sum_{s'} p(s'|s,a) \left[ E[r-r'|s,a,s'] + v^*(s') \right]$$

$r' = R_{t+k+1}$ is the end of episode reward

which we can't determine knowing $s, a, s'$

$\Rightarrow$ We can't write Bellman optimality eqn in this case

②

**SARSA :**

$$q_{new}(s,a) = q_{old}(s,a) + \alpha \left[ R_{t+1} + \gamma \, q_{old}(s',a') - q_{old}(s,a) \right]$$

Benefit : implicit evaluation and greedification , on-policy

Pitfall : explores too much - sometimes

**Expected SARSA :**

$$q(s_t, a_t) = q(s_t, a_t) + \alpha \left[ R_{t+1} + \gamma \sum_{a'} \pi(s_{t+1}, a') \, q(s_{t+1}, a') - q(s_t, a_t) \right]$$

Benefit : Converges to a more optimal q value

Pitfall : More Sampling required ✗ ?

$\alpha$ - learning

$$q(s_t, a_t) = q(s_t, a_t) + \alpha \left[ R_{t+1} + \gamma \max_{a'} q(s_{t+1}, a') - q(s_t, a_t) \right]$$

Benefit : Converges to the optimal solution more quickly ✗

Pitfall : Sometimes, if is too greedy

$$\eta(\theta) = E[R_t]$$

$$= \sum_a q_*(a) \, \pi(a, \theta)$$

$$Pr(a_t = a) = \frac{e^{q_t^*(a)/\beta}}{\sum_{b=1}^{n} e^{q_t(b)/\beta}} = \pi(a, \beta)$$

Thus, we have a definition for the policy $\pi$ parametorized by the temperature parameter $\beta$

The updates are

$$q_{t+1}(a_t) \doteq q_t(a_t) + \alpha(R_t - b_t)(1 - \pi_t(a_t)) \text{ and } \quad —①$$

$$q_{t+1}(a) \doteq q_t(a) - \alpha(R_t - b_t)\pi_t(a) \quad \forall a \neq a_t \quad —②$$

And the parameters can be updated as

$$\Delta\beta_t = \alpha_n(R_t - b_t)\frac{\partial \ln(a_t; \beta_t)}{\partial\theta}$$

Thus, we have casted this as a

⑥

policy gradient problem.

Yes, in some sense, $\gamma$ clarifies the eligibility

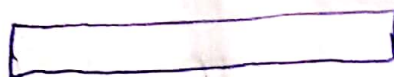of a particular reward term to the total return $G_t$

$$G_t = R_{t+1} + \gamma R_{t+2} + \boxed{\gamma^2 R_{t+3}} + \ldots$$

$\gamma$ gives indication of how much of $R_{t+3}$ is

valued (taken into account) into $G_t$. This ($\gamma$) gives some

kind of a special eligibility distribution, whereas $\lambda$
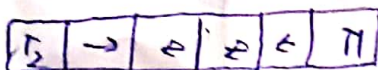
gives temporal eligibility.

⑧ There are total 16 policies possible, out
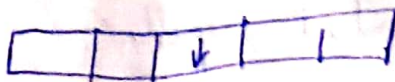
of which following are Blackwell optimal

①



$$V(s_1) = \frac{a\gamma}{1-\gamma^2} \qquad V(s_2) = \frac{a}{1-\gamma^2}, \qquad V(s_3) = \frac{a}{1-\gamma^2}$$

→ optimal for $a > 0$, $\gamma > \frac{-a + \sqrt{a^2 + 400}}{20}$
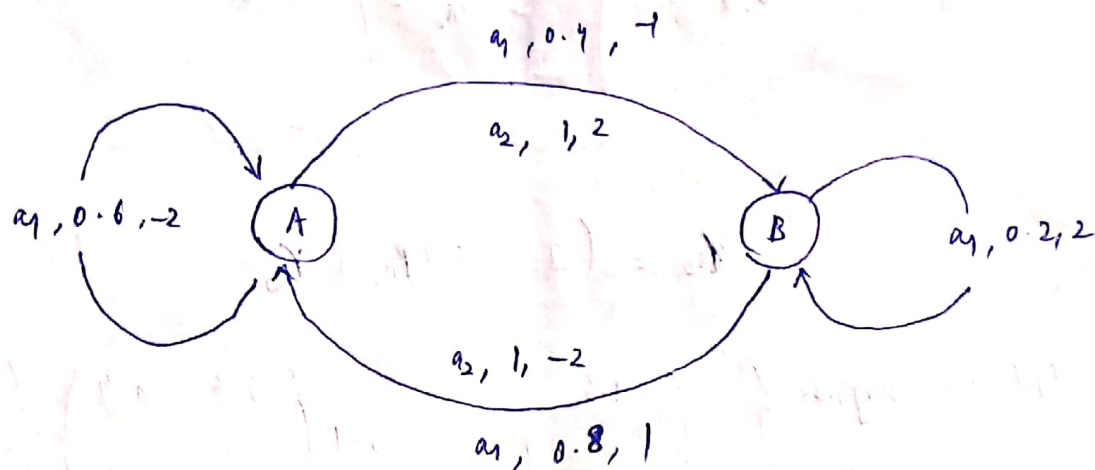
②



$$V(s) = 10\gamma^2$$
$$V(s) = 10\gamma^2$$
$$V(s) = 10$$
$$V(s) = 10$$

op fund for $\quad a < \frac{5}{\sqrt{2}}, \quad \boxed{?, > \frac{1}{\sqrt{2}} = k} \quad \boxed{k = 0.5}$



$a_1, 0.4, -1$

$a_2, 1, 2$

$a_1, 0.6, -2$    (A)        (B)    $a_1, 0.2, 2$

$a_2, 1, -2$

$a_1, 0.8, 1$

(a)     Initial Value $t : \{ A: 0, \quad B: 0 \} \qquad \gamma = 1$

$$V_1(A) = \max_a \left\{ a_1 : \overset{-1.6}{-8} + (1)(0.4)(0) + (1)(0.6)(-2, 0) \right\}$$

$a_2 :$

$$V_1(A) = \max_a \left\{ \begin{array}{l} a_1 : -1.6 + 1(0.4)(0) + (1)(0.6)(0) \\ a_2 : \quad 2 + (1)(1)(0) \end{array} \right\} = 2 \checkmark$$

$$V_1(B) = \max_a \left\{ \begin{array}{l} a_1 : +1.2 + (1)(0.2)(2) \\ \qquad\qquad + (1)(0.8)(2) \end{array} \right\} \checkmark$$

$$V_1(B) = \max \left\{ \begin{array}{l} a_1 : 1.2 + 0 + 0 \\ a_2 : -2 + 0 + 0 \end{array} \right\} = 1.2$$

$$V_2(A) = \max \left\{ \begin{array}{l} a_1 : -1.6 + 0.4\left(\frac{2}{0}\right) + 0.6\left(\frac{2}{0}\right) \\ a_2 : \quad 2 + 2 \end{array} \right\} = 4$$

$$V_2(B) = \max \left\{ \begin{array}{l} a_1 : -1.6 + 0.2(1.2) + 0.8(1.2) \\ a_2 : \quad 2 + (1)(1.2) \end{array} \right\} = 3.2$$

(b)

$$P_{\pi_0} = \begin{pmatrix} 0.6 & 0.4 \\ 0 & 1 \end{pmatrix} \qquad r_{\pi_0} = \begin{pmatrix} -1.6 \\ 2 \end{pmatrix}$$

$$V_{\pi_0} = (I - 0.9\, P_{\pi_0})\, r_{\pi_0}$$

$$\pi_1(A) = \arg\max_{a} \left\{ \begin{pmatrix} -1.6 \\ 2 \end{pmatrix} + 0.9 \begin{pmatrix} 0.6 & 0.4 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} 2 \\ 1.2 \end{pmatrix} \right\}$$

$$= \arg\max_{a} \left\{ \begin{pmatrix} -1.6 \\ 2 \end{pmatrix} + \begin{pmatrix} 1.51 \\ 1.33 \end{pmatrix} \right\} = \arg\max_{a} \left\{ \begin{matrix} -0.09 \\ 3.53 \end{matrix} \right\}$$

$$= a_2$$

$$\pi_2(B) = \arg\max_{a} \left\{ \begin{pmatrix} 1.2 \\ -2 \end{pmatrix} + 0.9 \begin{pmatrix} 0.2 & 0.8 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} \cdot \\ \cdot \end{pmatrix} \right\}$$

Sly, other iteration can be done

(1)

$$\langle \textit{s}, \qquad \langle n, f(n) \rangle$$

use errors as rewards