

Q1 ①

action & sequence: [1 0 0]

R: [-2 -8 9] \Rightarrow future return: [-1 1 9]

Model o/p

[0.59 0.41] P: 0.41 0.47 0.67

[0.47 0.53] log P: -0.892 -0.755 -0.4

[0.67 0.33]

$$\log P * R = 0.812 - 0.755 - 3.6$$

$$\text{Total loss} = - \left(\sum \right)$$

$$= |-3.463| = 3.463$$

② Action [0 0 0 2]

r [-8 -2 5 6]

P: 0.79 0.03 0.52 0.58

log P: -0.24 -3.51 -0.65 -0.58

\Rightarrow ①

$V_0 \quad V_1 \quad V_2 \quad V_3$

$$\Delta \text{ETA} = r + \delta u(s_t + 1) - V^{\pi} \quad -3.53 \quad -1.4 \quad -4.95 \quad -9.52$$

$$J_0 = -8 + (-1.4) + 3.53 = -5.87$$

$$J_1 = -2 + (-4.95) + 1.4 = -5.55 \quad J = [-5.87 \quad -5.55 \quad 0.43]$$

$$J_2 = 5 + (-9.52) + 4.95 = 0.43$$

Action loss [-1.41 -19.48 0.27]

$g = (9 * 5)$

Date 1 / 1

Subject _____

$$\text{critic loss} = [34.46 \quad 30.8 \quad 0.18]$$

\downarrow
 b^2

$$\text{overall loss} = [33.05 \quad 11.32 \quad 0.45]$$

\downarrow

actor loss + critic loss

$$\textcircled{3} \begin{array}{cccccc|cccc} a_0 & a_1 & a_2 & a_3 & a_4 & r_0 & r_1 & r_2 & r_3 & r_4 \\ 3 & 2 & 3 & 4 & 3 & 5 & -1 & -3 & 4 & 8 \end{array}$$

$$\begin{array}{ccccc|ccccc} T(0) & T(1) & T(2) & T(3) & T(4) & v_0 & v_1 & v_2 & v_3 & v_4 \\ -5.86 & -3.98 & -2.89 & -7 & -2.89 & -3.88 & -2.17 & -9.26 & -4.72 & -1.73 \end{array}$$

$$\delta = 0.97, N = 0.99, E = 0.2$$

① Advantage:

$$G = r + \delta \cdot v(s+i) - v(s+1)$$

$$G_0 = 5 + 0.97 \times -2.17 + 3.88 = 6.78$$

$$G_1 = -1 + 0.97 \times -9.26 + 2.17 = -7.81$$

$$G_2 = -3 + 0.97 \times -4.72 + 9.26 = 1.682$$

$$G_3 = 4 + 0.97 \times -1.73 + 4.72 = 10.4$$

$$A = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 1 & 0 & 0 & 1 \end{bmatrix} \begin{matrix} \delta N \\ \delta^2 N^2 \\ \delta^3 N^3 \\ \delta^4 N^4 \end{matrix} * \begin{bmatrix} G_0 \\ G_1 \\ G_2 \\ G_3 \end{bmatrix} = \begin{bmatrix} 6.78 \\ -7.81 \\ 1.682 \\ 10.4 \end{bmatrix}$$

$$A_1 = 6.78 + (0.97 \times 0.99 \times -7.81) + (0.97^2 \times 0.99^2 \times 1.682) + (0.97^3 \times 0.99^3 \times 10.4) = 10.04$$

$$A_2 = -7.81 + (0.97 \times 0.99 \times 1.682) + (0.97^2 \times 0.99^2 \times 10.4) = 3.4$$

$$A_3 = 1.682 + (0.97 \times 0.99 \times 10.4) = 11.67$$

$$A_4 = 10.4$$

A_0	A_1	A_2	A_3
10.04	3.4	11.67	10.4

② actor loss

P_{Told}	0.39	0.07	0.06	0.18	0.04
P_{Tnew}	0.17	0.14	0.17	0.36	0.24
$Ration \frac{new}{old}$					

RT:

R	r_1	r_2	r_3	r_4
0.438	2	2.83	2	6

$$C = 0.2 \quad \begin{cases} 1 - C = 0.8 \\ 1 + C = 1.2 \end{cases}$$

Clipped RT:

r_{ro}	r_1	r_2	r_3	r_4
0.8	1.2	1.2	1.2	1.2

RT AT:	0	1	2	3	4
	4.38	6.8	33.03	20.8	0

RT Clipped AT:	0	1	2	3	4
	6.03	4.03	14	12.48	0

Min (clipped - unclipped)

0	1	2	3	4
4.38	4.08	14	12.48	0

$$\text{actor loss} = \text{mean} = 6.982$$

③ critic loss

Return: $P + V_{old}$

0	1	2	3	4
4.18	-0.58	3.78	3.4	-2.89

Return - V_{new}

0	1	2	3	4
8.06	1.59	8.04	8.12	-1.16

Spd Return - V_{new}

0	1	2	3	4
64.96	7.53	170.04	65.93	1.35

critical loss = Mean = 60.96