

111550177_HW2

1.

晚上 8:34 4月19日 週五

111550177_RL_HW2

Problem 1 (Baseline for Variance Reduction) (8+8+8=24 points)

Consider an example similar to that in the slides of Lecture 8 for explaining the baseline. Suppose there are only 1 non-terminal starting state (denoted by s) and 3 actions (denoted by a, b, c) in the MDP of interest. After any one of the action is applied at the starting state s , the MDP would evolve from s to the terminal state, with probability 1. Moreover, consider the following setting:

- The rewards are deterministic, and the reward function is defined as $r(s, a) = 100$, $r(s, b) = 98$, and $r(s, c) = 95$. Moreover, there is no terminal reward.
- We consider a softmax policy with parameters $\theta_a, \theta_b, \theta_c$ such that $\pi_\theta(\cdot|s) = \exp(\theta)/(\exp(\theta_a) + \exp(\theta_b) + \exp(\theta_c))$. Moreover, currently the parameters are $\theta_a = 0$, $\theta_b = \ln 5$, $\theta_c = \ln 4$.
- We would like to combine PG with SGD. At each policy update, we would construct an unbiased estimate $\hat{\nabla}V$ of the true policy gradient $\nabla_\theta V^{\pi_\theta}$ by sampling one trajectory (Note: $\hat{\nabla}V$ is a random vector. In this example, each trajectory has only one time step, and $s_0 = s$, a_0 is either a, b , or c , and s_1 is the terminal state).

(a) What are the mean vector of $\hat{\nabla}V$ (denoted by $\mathbb{E}[\hat{\nabla}V]$) and the covariance matrix of $\hat{\nabla}V$ (i.e., $\mathbb{E}[(\hat{\nabla}V - \mathbb{E}[\hat{\nabla}V])(\hat{\nabla}V - \mathbb{E}[\hat{\nabla}V])^\top]$)?

(a)

$$\pi_\theta(a|s) = \frac{e^{\theta_a}}{e^{\theta_a} + e^{\theta_b} + e^{\theta_c}} = \frac{e^0}{e^0 + e^{(\ln 5)} + e^{(\ln 4)}} = \frac{1}{10}$$

$$\pi_\theta(b|s) = \frac{e^{\theta_b}}{e^{\theta_a} + e^{\theta_b} + e^{\theta_c}} = \frac{e^{(\ln 5)}}{e^0 + e^{(\ln 5)} + e^{(\ln 4)}} = \frac{5}{10}$$

$$\pi_\theta(c|s) = \frac{e^{\theta_c}}{e^{\theta_a} + e^{\theta_b} + e^{\theta_c}} = \frac{e^{(\ln 4)}}{e^0 + e^{(\ln 5)} + e^{(\ln 4)}} = \frac{4}{10}$$

$$\nabla \log \pi_\theta(a|s) = \begin{pmatrix} \frac{\partial \log \pi_\theta(a|s)}{\partial \theta_a} \\ \frac{\partial \log \pi_\theta(a|s)}{\partial \theta_b} \\ \frac{\partial \log \pi_\theta(a|s)}{\partial \theta_c} \end{pmatrix} = \begin{pmatrix} 1 - \pi_\theta(a|s) \\ -\pi_\theta(b|s) \\ -\pi_\theta(c|s) \end{pmatrix} = \begin{pmatrix} 1 - 0.1 \\ -0.5 \\ -0.4 \end{pmatrix} = \begin{pmatrix} 0.9 \\ -0.5 \\ -0.4 \end{pmatrix}$$

$$\nabla \log \pi_\theta(b|s) = \begin{pmatrix} \frac{\partial \log \pi_\theta(b|s)}{\partial \theta_a} \\ \frac{\partial \log \pi_\theta(b|s)}{\partial \theta_b} \\ \frac{\partial \log \pi_\theta(b|s)}{\partial \theta_c} \end{pmatrix} = \begin{pmatrix} -\pi_\theta(a|s) \\ 1 - \pi_\theta(b|s) \\ -\pi_\theta(c|s) \end{pmatrix} = \begin{pmatrix} -0.1 \\ 1 - 0.5 \\ -0.4 \end{pmatrix} = \begin{pmatrix} -0.1 \\ 0.5 \\ -0.4 \end{pmatrix}$$

111550177_RL_HW2

數值方法_assi... 電子表格 Ch2 電子表格 Ch3 數值方法 HW2 111550177...

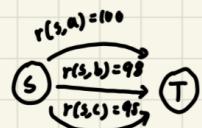


$$\nabla \log \pi_\theta(a|s) = \begin{pmatrix} \frac{\partial \log \pi_\theta(a|s)}{\partial \theta_a} \\ \frac{\partial \log \pi_\theta(a|s)}{\partial \theta_b} \\ \frac{\partial \log \pi_\theta(a|s)}{\partial \theta_c} \end{pmatrix} = \begin{pmatrix} -\pi_\theta(a|s) \\ -\pi_\theta(b|s) \\ 1-\pi_\theta(c|s) \end{pmatrix} = \begin{pmatrix} -0.1 \\ -0.5 \\ 1-0.4 \end{pmatrix} = \begin{pmatrix} -0.1 \\ -0.5 \\ 0.6 \end{pmatrix}$$

$$E[\hat{v}] = \nabla_\theta V^\pi$$

$$= E_{\tau \sim P_\theta^{\pi_0}} \left[\sum_{t=0}^{\infty} \gamma^t Q^{\pi_0}(s_t, a_t) \nabla_\theta \log \pi_\theta(a_t|s_t) \right]$$

$$= \sum_{a_0 \in \{a, b, c\}} r(s_0, a_0) \cdot \nabla_\theta \log \pi_\theta(a_0|s) \cdot \pi_\theta(a_0|s)$$



$$= 0.1 \times 100 \times \begin{pmatrix} 0.9 \\ -0.5 \\ -0.4 \end{pmatrix} + 0.5 \times 98 \times \begin{pmatrix} -0.1 \\ 0.5 \\ -0.4 \end{pmatrix} + 0.4 \times 95 \times \begin{pmatrix} -0.1 \\ -0.5 \\ 0.6 \end{pmatrix}$$

$$= \begin{pmatrix} 0.3 \\ 0.5 \\ -0.8 \end{pmatrix}$$

$$E[(\hat{v} - E[\hat{v}]) \cdot (\hat{v} - E[\hat{v}])^T]$$

$$= E_{\tau \sim P_\theta^{\pi_0}} [(\hat{v} - E[\hat{v}]) \cdot (\hat{v} - E[\hat{v}])^T]$$

$$= \sum_{a_0 \in \{a, b, c\}} (\hat{v}_{a_0} - E[\hat{v}]) (\hat{v}_{a_0} - E[\hat{v}])^T \cdot \pi(a_0|s)$$

$$= \left(100 \times \begin{pmatrix} 0.9 \\ -0.5 \\ -0.4 \end{pmatrix} - \begin{pmatrix} 0.3 \\ 0.5 \\ -0.8 \end{pmatrix} \right) \cdot \left(100 \times \begin{pmatrix} 0.9 \\ -0.5 \\ -0.4 \end{pmatrix} - \begin{pmatrix} 0.3 \\ 0.5 \\ -0.8 \end{pmatrix} \right)^T \times 0.1$$

晚上 8:35 4月19日 週五

111550177_RL_HW2

數值方法_assi... 雖值方法 Ch2 雖值方法 Ch3 數值方法 HW2 111550177...

$$+ (q_8 \times \begin{bmatrix} -0.1 \\ 0.5 \\ -0.4 \end{bmatrix} - \begin{bmatrix} 0.3 \\ 0.5 \\ -0.8 \end{bmatrix}) \cdot (q_8 \times \begin{bmatrix} -0.1 \\ 0.5 \\ -0.4 \end{bmatrix} - \begin{bmatrix} 0.3 \\ 0.5 \\ -0.8 \end{bmatrix})^T \times 0.5$$

$$+ (q_5 \times \begin{bmatrix} -0.1 \\ -0.5 \\ 0.6 \end{bmatrix} - \begin{bmatrix} 0.3 \\ 0.5 \\ -0.8 \end{bmatrix}) \cdot (q_5 \times \begin{bmatrix} -0.1 \\ -0.5 \\ 0.6 \end{bmatrix} - \begin{bmatrix} 0.3 \\ 0.5 \\ -0.8 \end{bmatrix})^T \times 0.4$$

$$= \begin{bmatrix} 894.03 & -509.75 & -384.28 \\ -509.75 & 2352.75 & -1834 \\ -384.28 & -1834 & 2227.28 \end{bmatrix}$$

(b) Suppose we leverage the value function $V^{\pi_\theta}(s)$ as the baseline and denote by $\tilde{\nabla}V$ the corresponding estimated policy gradient. Then, what are the mean vector and the covariance matrix of $\tilde{\nabla}V$? (Note: $\tilde{\nabla}V$ is also a random vector)

$$V^{\pi_\theta}(s) = 0.1 \times 100 + 0.5 \times 98 + 0.4 \times 95 = 97$$

$$\mathbb{E}[\tilde{\nabla}V] = \nabla_{\theta} V^{\pi_\theta}$$

$$= \mathbb{E}_{\tau \sim P_{\theta}^{\pi_\theta}} \left[\sum_{t=0}^{\infty} r_t \left(Q^{\pi_\theta}(s_t, a_t) - V^{\pi_\theta}(s_t) \right) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \right]$$

$$= \sum_{a_t \in \{a_0, a_1, a_2\}} (r(s, a_t) - V^{\pi_\theta}) \cdot \nabla_{\theta} \log \pi_{\theta}(a_t | s) \cdot \pi_{\theta}(a_t | s)$$

$$= 0.1 \times (100 - 97) \times \begin{bmatrix} 0.9 \\ -0.5 \\ -0.4 \end{bmatrix} + 0.5 \times (98 - 97) \times \begin{bmatrix} -0.1 \\ 0.5 \\ -0.4 \end{bmatrix} + 0.4 \times (95 - 97) \times \begin{bmatrix} -0.1 \\ -0.5 \\ 0.6 \end{bmatrix}$$

第3頁，共9頁



$$= \begin{pmatrix} 0.3 \\ 0.5 \\ -0.8 \end{pmatrix}$$

$$E[(\tilde{v} - E[\tilde{v}]) \cdot (\tilde{v} - E[\tilde{v}])^T]$$

$$= E_{v \sim P_{a_0}^{x_0}} [(\tilde{v} - E[\tilde{v}]) \cdot (\tilde{v} - E[\tilde{v}])^T]$$

$$= \sum_{a_0 \in \{a_0, b_0\}} (\tilde{v}_{a_0} - E[\tilde{v}]) (\tilde{v}_{a_0} - E[\tilde{v}])^T \cdot \pi(a_0 | s)$$

$$= ((100-99) \times \begin{pmatrix} 0.9 \\ -0.5 \\ 0.5 \\ -0.4 \end{pmatrix} - \begin{pmatrix} 0.3 \\ 0.5 \\ -0.8 \end{pmatrix}) \cdot ((100-99) \times \begin{pmatrix} 0.9 \\ -0.5 \\ 0.5 \\ -0.4 \end{pmatrix} - \begin{pmatrix} 0.3 \\ 0.5 \\ -0.8 \end{pmatrix})^T \times 0.1$$

$$+ ((98-99) \times \begin{pmatrix} -0.1 \\ 0.5 \\ -0.4 \end{pmatrix} - \begin{pmatrix} 0.3 \\ 0.5 \\ -0.8 \end{pmatrix}) \cdot ((98-99) \times \begin{pmatrix} -0.1 \\ 0.5 \\ -0.4 \end{pmatrix} - \begin{pmatrix} 0.3 \\ 0.5 \\ -0.8 \end{pmatrix})^T \times 0.5$$

$$+ ((95-99) \times \begin{pmatrix} -0.1 \\ -0.5 \\ 0.6 \end{pmatrix} - \begin{pmatrix} 0.3 \\ 0.5 \\ -0.8 \end{pmatrix}) \cdot ((95-99) \times \begin{pmatrix} -0.1 \\ -0.5 \\ 0.6 \end{pmatrix} - \begin{pmatrix} 0.3 \\ 0.5 \\ -0.8 \end{pmatrix})^T \times 0.4$$

$$= \begin{pmatrix} 0.66 & -0.5 & -0.16 \\ -0.5 & 0.5 & 0 \\ -0.16 & 0 & 0.16 \end{pmatrix}$$

第4頁，共9頁



(c) Let $B(s)$ denote a baseline function and ∇V_B be the corresponding estimated policy gradient (∇V_B is again a random vector). Suppose we say that a baseline function $B(s)$ is *optimal* if it attains the minimum trace of the corresponding covariance matrix of ∇V_B among all possible state-dependent baselines. Please try to find one such optimal $B(s)$.

The covariance matrix with baseline is

$$E[\tilde{v}] = 0.1 \times (100 - B) \times \begin{pmatrix} 0.9 \\ -0.5 \\ -0.4 \end{pmatrix} + 0.5 \times (98 - B) \times \begin{pmatrix} -0.1 \\ 0.5 \\ -0.4 \end{pmatrix} + 0.4 \times (95 - B) \times \begin{pmatrix} -0.1 \\ -0.5 \\ 0.6 \end{pmatrix} = \begin{pmatrix} 0.3 \\ 0.5 \\ -0.8 \end{pmatrix}$$

$$E[(\tilde{v} - E[\tilde{v}]) \cdot (\tilde{v} - E[\tilde{v}])^T] = \begin{pmatrix} a & b & c \\ b & b^2 & bc \\ c & bc & c^2 \end{pmatrix}$$

$$= ((100 - B) \times \begin{pmatrix} 0.9 \\ -0.5 \\ -0.4 \end{pmatrix} - \begin{pmatrix} 0.3 \\ 0.5 \\ -0.8 \end{pmatrix}) \cdot ((100 - B) \times \begin{pmatrix} 0.9 \\ -0.5 \\ -0.4 \end{pmatrix} - \begin{pmatrix} 0.3 \\ 0.5 \\ -0.8 \end{pmatrix})^T \times 0.1$$

$$+ ((98 - B) \times \begin{pmatrix} -0.1 \\ 0.5 \\ -0.4 \end{pmatrix} - \begin{pmatrix} 0.3 \\ 0.5 \\ -0.8 \end{pmatrix}) \cdot ((98 - B) \times \begin{pmatrix} -0.1 \\ 0.5 \\ -0.4 \end{pmatrix} - \begin{pmatrix} 0.3 \\ 0.5 \\ -0.8 \end{pmatrix})^T \times 0.5$$

$$+ ((95 - B) \times \begin{pmatrix} -0.1 \\ -0.5 \\ 0.6 \end{pmatrix} - \begin{pmatrix} 0.3 \\ 0.5 \\ -0.8 \end{pmatrix}) \cdot ((95 - B) \times \begin{pmatrix} -0.1 \\ -0.5 \\ 0.6 \end{pmatrix} - \begin{pmatrix} 0.3 \\ 0.5 \\ -0.8 \end{pmatrix})^T \times 0.4$$

The trace of covariance matrix :

$$\text{tr}(B(s)) = 0.1 \times \left\{ [(100 - B(s)) \times 0.9 - 0.3]^2 + [(100 - B(s)) \times (-0.5) - 0.5]^2 + [(100 - B(s)) \times (-0.4) + 0.8]^2 \right\}$$

晚上 8:35 4月19日 週五

111550177_RL_HW2

$+ 0.5 \times \left[(98 - B(s)) \times (-0.1) - 0.3 \right]^2 + \left[(98 - B(s)) \times 0.5 - 0.5 \right]^2 + \left[(98 - B(s)) \times (-0.4) + 0.8 \right]^2 \right]$
 $+ 0.4 \left\{ \left[(95 - B(s)) \times (-0.1) - 0.3 \right]^2 + \left[(95 - B(s)) \times (-0.5) - 0.5 \right]^2 + \left[(95 - B(s)) \times 0.6 + 0.8 \right]^2 \right\}$

 $\frac{df(B(s))}{dB(s)} = 0 \Rightarrow 0.2 \times (121.66 - 1.22B(s)) + 1 \times (40.62 - 0.42 \cdot B(s)) + 0.8 \times (59.66 - 0.62B(s)) = 0$
 $\Rightarrow B(s) = 97.139931 \text{ is the optimal Baseline}$

Problem 2 (Non-Uniform Polyak-Lojacsiewicz Condition in RL) (8+8=16 points)

As described in Lecture 12, let us prove the fundamental Polyak-Lojacsiewicz condition in RL: Let π^* be an optimal policy and let $a^* := \arg \max_a \pi^*(a|s)$ (essentially, a^* is an optimal action). Under softmax policies, we have

$$\left\| \frac{\partial V^{\pi_\theta}(\mu)}{\partial \theta} \right\|_2 \geq \frac{1}{\sqrt{S}} \left\| \frac{d\pi^*}{d\mu} \right\|_\infty^{-1} \cdot \min_{s \in S} \pi_\theta(a^*(s)|s) \cdot [V^*(\mu) - V^{\pi_\theta}(\mu)]. \quad (1)$$

To show this, you would also need the celebrated “Performance difference lemma” as follows: For any two policies π_1 and π_2 , we always have

$$V^{\pi_2}(\mu) - V^{\pi_1}(\mu) = \frac{1}{1-\gamma} \mathbb{E}_{s' \sim d_\mu^{\pi_2}} \mathbb{E}_{a' \sim \pi_2(\cdot|s')} [A^{\pi_1}(s', a')]. \quad (2)$$

(a) To begin with, show the following result:

$$\left\| \frac{\partial V^{\pi_\theta}(\mu)}{\partial \theta} \right\|_2 \geq \frac{1}{1-\gamma} \cdot \frac{1}{\sqrt{S}} \sum_s d_\mu^{\pi_\theta}(s) \cdot \pi_\theta(a^*(s)|s) \cdot |A^{\pi_\theta}(s, a^*(s))|. \quad (3)$$

(Hint: You would need to first apply Cauchy-Schwarz inequality and leverage the Policy Gradient expression under softmax policies. This subproblem shall require about 5-8 lines of proof.)

$$\begin{aligned} \left\| \frac{\partial V^{\pi_\theta}(\mu)}{\partial \theta} \right\|_2 &\geq \left[\sum_{s \in S, a \in A} \left(\frac{\partial V^{\pi_\theta}(\mu)}{\partial \theta(s, a)} \right)^2 \right]^{\frac{1}{2}} \\ &\geq \left[\sum_{s \in S} \left(\frac{\partial V^{\pi_\theta}(\mu)}{\partial \theta(s, a^*)} \right)^2 \right]^{\frac{1}{2}} \\ &\geq \frac{1}{\sqrt{S}} \sum_{s \in S} \left| \frac{\partial V^{\pi_\theta}(\mu)}{\partial \theta(s, a^*)} \right| \end{aligned}$$

For a vector $Y = (y_1, y_2, y_3, \dots, y_n)$, $\|Y\|_1 = \sum_{i=1}^n |y_i|$ is a vector $(\underbrace{1, 1, \dots, 1})$
 $\|Y\|_2 = \sqrt{\sum_{i=1}^n y_i^2} = \sqrt{\langle Y, 1 \rangle}$ by Cauchy-Schwarz inequality

第6頁，共9頁 PG under softmax policies

2.

晚上 8:36 4月19日 週五 100% 111550177_RL_HW2

數值方法_assi... × 數值方法 Ch2 × 數值方法 Ch3 × 數值方法 HW2 × 111550177...

Problem 2 (Non-Uniform Polyak-Lojacsiewicz Condition in RL) (8+8=16 points)

As described in Lecture 12, let us prove the fundamental Polyak-Lojacsiewicz condition in RL: Let π^* be an optimal policy and let $a^* := \arg \max_a \pi^*(a|s)$ (essentially, a^* is an optimal action). Under softmax policies, we have

$$\left\| \frac{\partial V^{\pi_\theta}(\mu)}{\partial \theta} \right\|_2 \geq \frac{1}{\sqrt{S}} \left\| \frac{d_\mu^{\pi^*}}{d_\mu^{\pi_\theta}} \right\|_\infty^{-1} \cdot \min_{s \in S} \pi_\theta(a^*(s)|s) \cdot [V^*(\mu) - V^{\pi_\theta}(\mu)]. \quad (1)$$

To show this, you would also need the celebrated “Performance difference lemma” as follows: For any two policies π_1 and π_2 , we always have

$$V^{\pi_2}(\mu) - V^{\pi_1}(\mu) = \frac{1}{1-\gamma} \mathbb{E}_{s' \sim d_\mu^{\pi_2}} \mathbb{E}_{a' \sim \pi_2(\cdot|s')} [A^{\pi_1}(s', a')]. \quad (2)$$

(a) To begin with, show the following result:

$$\left\| \frac{\partial V^{\pi_\theta}(\mu)}{\partial \theta} \right\|_2 \geq \frac{1}{1-\gamma} \cdot \frac{1}{\sqrt{S}} \sum_s d_\mu^{\pi_\theta}(s) \cdot \pi_\theta(a^*(s)|s) \cdot |A^{\pi_\theta}(s, a^*(s))|. \quad (3)$$

(Hint: You would need to first apply Cauchy-Schwarz inequality and leverage the Policy Gradient expression under softmax policies. This subproblem shall require about 5-8 lines of proof.)

$$\begin{aligned} \left\| \frac{\partial V^{\pi_\theta}(\mu)}{\partial \theta} \right\|_2 &\geq \left[\sum_{s \in S, a \in A} \left(\frac{\partial V^{\pi_\theta}(\mu)}{\partial \theta(s, a)} \right)^2 \right]^{\frac{1}{2}} \\ &\geq \left[\sum_{s \in S} \left(\frac{\partial V^{\pi_\theta}(\mu)}{\partial \theta(s, a^*)} \right)^2 \right]^{\frac{1}{2}} \\ &\geq \frac{1}{\sqrt{S}} \sum_{s \in S} \left| \frac{\partial V^{\pi_\theta}(\mu)}{\partial \theta(s, a^*)} \right| \end{aligned}$$

For a vector $Y = (y_1, y_2, y_3, \dots, y_n)$, $\|Y\|_1 = \sum_{i=1}^n y_i$ is a vector $\underbrace{(1, 1, \dots, 1)}_{n \times 1}$

$\|Y\|_1 = \sum_{i=1}^n y_i = \sum_{i=1}^n y_i \cdot 1 = |\langle |Y|, 1 \rangle|$ by Cauchy-Schwarz inequality

$\leq \|Y\|_2 \cdot \|1\|$

PG under softmax policies

$$\begin{aligned} &= \frac{1}{1-\gamma} \cdot \frac{1}{\sqrt{S}} \sum_{s \in S} |d_\mu^{\pi_\theta}(s) \cdot \pi_\theta(a^*(s)|s) \cdot A^{\pi_\theta}(s, a^*(s))| \\ &= \frac{1}{1-\gamma} \cdot \frac{1}{\sqrt{S}} \sum_{s \in S} d_\mu^{\pi_\theta}(s) \cdot \pi_\theta(a^*(s)|s) \cdot |A^{\pi_\theta}(s, a^*(s))| \quad (d_\mu^{\pi_\theta}(s), \pi_\theta(a^*(s)|s) \text{ is non-negative}) \end{aligned}$$

(b) Next, please use the results in (a) and the Performance difference lemma to conclude that the PL condition in (1) indeed holds. (Hint: Try to handle each term in (3) separately. This subproblem shall require only about 5-8 lines of proof.)

$$\begin{aligned}
 &= \frac{1}{1-\gamma} \cdot \frac{1}{\sqrt{s}} \sum_{s \in S} |d_{\mu}^{k_0}(s) \cdot \pi_{\theta}(\alpha^*(s)|s) \cdot A^{k_0}(s, \alpha^*(s))| \\
 &= \frac{1}{1-\gamma} \cdot \frac{1}{\sqrt{s}} \sum_{s \in S} d_{\mu}^{k_0}(s) \cdot \pi_{\theta}(\alpha^*(s)|s) \cdot |A^{k_0}(s, \alpha^*(s))| \quad (d_{\mu}^{k_0}(s), \pi_{\theta}(\alpha^*(s)|s) \\
 &\quad \text{is non-negative})
 \end{aligned}$$

(b) Next, please use the results in (a) and the Performance difference lemma to conclude that the PL condition in (1) indeed holds. (Hint: Try to handle each term in (3) separately. This subproblem shall require only about 5-8 lines of proof.)

$$\begin{aligned}
 \left\| \frac{\partial V^{k_0}(\mu)}{\partial \theta} \right\|_2 &\geq \frac{1}{1-\gamma} \cdot \frac{1}{\sqrt{s}} \sum_{s \in S} d_{\mu}^{k_0}(s) \cdot \pi_{\theta}(\alpha^*(s)|s) \cdot |A^{k_0}(s, \alpha^*(s))| \\
 \left\| \frac{d_{\rho}^{k^*}(s)}{d_{\mu}^{k_0}(s)} \right\|_{\infty} &\geq \frac{1}{1-\gamma} \cdot \frac{1}{\sqrt{s}} \sum_{s \in S} \frac{d_{\mu}^{k_0}(s)}{d_{\mu}^{k_0}(s)} \cdot d_{\rho}^{k^*}(s) \cdot \pi_{\theta}(\alpha^*(s)|s) \cdot |A^{k_0}(s, \alpha^*(s))| \\
 = \max_s \frac{d_{\rho}^{k^*}(s)}{d_{\mu}^{k_0}(s)} &\geq \frac{1}{1-\gamma} \cdot \frac{1}{\sqrt{s}} \left\| \frac{d_{\rho}^{k^*}(s)}{d_{\mu}^{k_0}(s)} \right\|_{\infty}^{-1} \cdot \min_s \pi_{\theta}(\alpha^*(s)|s) \cdot \sum_{s \in S} d_{\rho}^{k^*}(s) \cdot |A^{k_0}(s, \alpha^*(s))| \\
 &\geq \frac{1}{1-\gamma} \cdot \frac{1}{\sqrt{s}} \left\| \frac{d_{\rho}^{k^*}(s)}{d_{\mu}^{k_0}(s)} \right\|_{\infty}^{-1} \cdot \min_s \pi_{\theta}(\alpha^*(s)|s) \cdot \sum_{s \in S} d_{\rho}^{k^*}(s) \cdot |A^{k_0}(s, \alpha^*(s))| \quad \text{allow negative} \\
 &= \frac{1}{\sqrt{s}} \left\| \frac{d_{\rho}^{k^*}(s)}{d_{\mu}^{k_0}(s)} \right\|_{\infty}^{-1} \cdot \min_s \pi_{\theta}(\alpha^*(s)|s) \cdot \frac{1}{1-\gamma} \cdot \sum_{s \in S} d_{\rho}^{k^*}(s) \cdot \underbrace{\sum_{a \in A} \pi_a(s) \cdot A^{k_0}(s, a)}_{\sim 1} \\
 &= \frac{1}{\sqrt{s}} \left\| \frac{d_{\rho}^{k^*}(s)}{d_{\mu}^{k_0}(s)} \right\|_{\infty}^{-1} \cdot \min_s \pi_{\theta}(\alpha^*(s)|s) \cdot \frac{1}{1-\gamma} \cdot E_{s \sim d_{\rho}^{k^*}} [E_{a \sim \pi_{\theta}} [A^{k_0}(s, a)]] \\
 &= \frac{1}{\sqrt{s}} \left\| \frac{d_{\rho}^{k^*}(s)}{d_{\mu}^{k_0}(s)} \right\|_{\infty}^{-1} \cdot \min_s \pi_{\theta}(\alpha^*(s)|s) \cdot \frac{1}{1-\gamma} \cdot [V^*(\mu) - V^{k_0}(\mu)] \quad \text{Performance Difference Lemma}
 \end{aligned}$$

晚上 8:37 4月19日 週五

100%

111550177_RL_HW2

Problem 3 (Monte Carlo Policy Evaluation) (8+8=16 points)

As discussed in Lecture 10, we know that the every-visit Monte Carlo estimate is biased. Let us quickly verify this fact under the simple 2-state MRP (which serves as a reduction from any MRP), as shown below. Specifically, we need to check the following two properties:

Figure 1: A simple 2-state MRP.

- Property 1: Show that the true value function at state S (denoted by $V(S)$) satisfies that

$$V(S) = \frac{P_S}{P_T} R_S + R_T. \quad (4)$$

Possible scenario :

$$S \rightarrow T : P_T R_T$$

$$S \rightarrow S \rightarrow T : P_S P_T (R_S + R_T)$$

$$S \rightarrow S \rightarrow S \rightarrow T : P_S^2 P_T (2R_S + R_T)$$

$$\vdots$$

$$S \rightarrow S \rightarrow \dots \rightarrow S \rightarrow T : P_S^n P_T (n \cdot R_S + R_T)$$

$$\Rightarrow V(S) = P_T R_T + P_S P_T (R_S + R_T) + P_S^2 P_T (2R_S + R_T) + \dots + P_S^n P_T (n \cdot R_S + R_T)$$

$$= \sum_{n=0}^{\infty} P_S^n P_T (n \cdot R_S + R_T)$$

$$= P_T \cdot R_S \cdot \sum_{n=0}^{\infty} n \cdot P_S^n + P_T \cdot R_T \cdot \sum_{n=0}^{\infty} P_S^n$$

$$= P_T \cdot R_S \cdot \frac{P_S}{(1-P_S)^2} + P_T \cdot R_T \cdot \frac{1}{1-P_S}$$

$\sum_{i=0}^{\infty} a^i = \frac{1}{1-a}, \sum_{i=0}^{\infty} i \cdot a^i = \frac{a}{(1-a)^2}, \text{ if } |a| < 1$
 $(P_S + P_T = 1)$

晚上 8:37 4月19日 週五

111550177_RL_HW2

$= \frac{P_S}{P_T} \cdot R_S + R_T$

- Property 2: Suppose we construct an every-visit MC estimate based on only 1 trajectory τ (denoted by $\hat{V}_{MC}(S; \tau)$). Then, please show that

$$\mathbb{E}_\tau[\hat{V}_{MC}(S; \tau)] = \frac{P_S}{2P_T}R_S + R_T. \quad (5)$$

(Hint: To begin with, you shall consider all possible trajectories and the corresponding probabilities. Accordingly, you would obtain that $\mathbb{E}_\tau[\hat{V}_{MC}(S; \tau)] = \sum_{k=0}^{\infty} P_T P_S^k \left(\frac{R_S + 2R_S + \dots + kR_S + (k+1)R_T}{k+1} \right).$)

Possible scenario:

$S \rightarrow T : P_T R_T$

$S \rightarrow S \rightarrow T : P_T P_S \cdot \left(\frac{R_S + 2R_T}{2} \right)$

$S \rightarrow S \rightarrow S \rightarrow T : P_T P_S^2 \cdot \left(\frac{R_S + 2R_S + 3R_T}{3} \right)$

⋮

$S \rightarrow S \rightarrow \dots \rightarrow S \rightarrow T : P_T P_S^n \cdot \left(\frac{R_S + 2R_S + \dots + nR_S + (n+1)R_T}{n+1} \right)$

$$\mathbb{E}_\tau[\hat{V}_{MC}(S; \tau)] = \sum_{n=0}^{\infty} P_T P_S^n \left(\frac{R_S + 2R_S + \dots + nR_S + (n+1)R_T}{n+1} \right)$$

$$= \sum_{n=0}^{\infty} P_T P_S^n \cdot \left[\frac{\frac{n(n+1)}{2}}{n+1} + R_T \right]$$

$$= \frac{P_T}{2} \sum_{n=0}^{\infty} P_S^n \cdot n + P_T R_T \sum_{n=0}^{\infty} P_S^n$$

$$= \frac{P_T}{2} \cdot \frac{P_S}{(1-P_S)^2} + P_T \cdot R_T \cdot \frac{1}{1-P_S}$$

$$= \frac{P_S}{2P_T} R_S + R_T$$

$\sum_{i=0}^{\infty} a^i = \frac{1}{1-a}, \sum_{i=0}^{\infty} i \cdot a^i = \frac{a}{(1-a)^2}$
, if $|a| < 1$
 $(P_S + P_T = 1)$

4.

(a)

I add two shared layer in my neuro network, and the activation function are both ReLu. I also add one dropout between layer one and two with dropout probability 0.2.

I found that when the gamma goes down, the episode needed to converge become more. The best number of hidden layer is 256. If the learning rate is too large, the model can't converge.

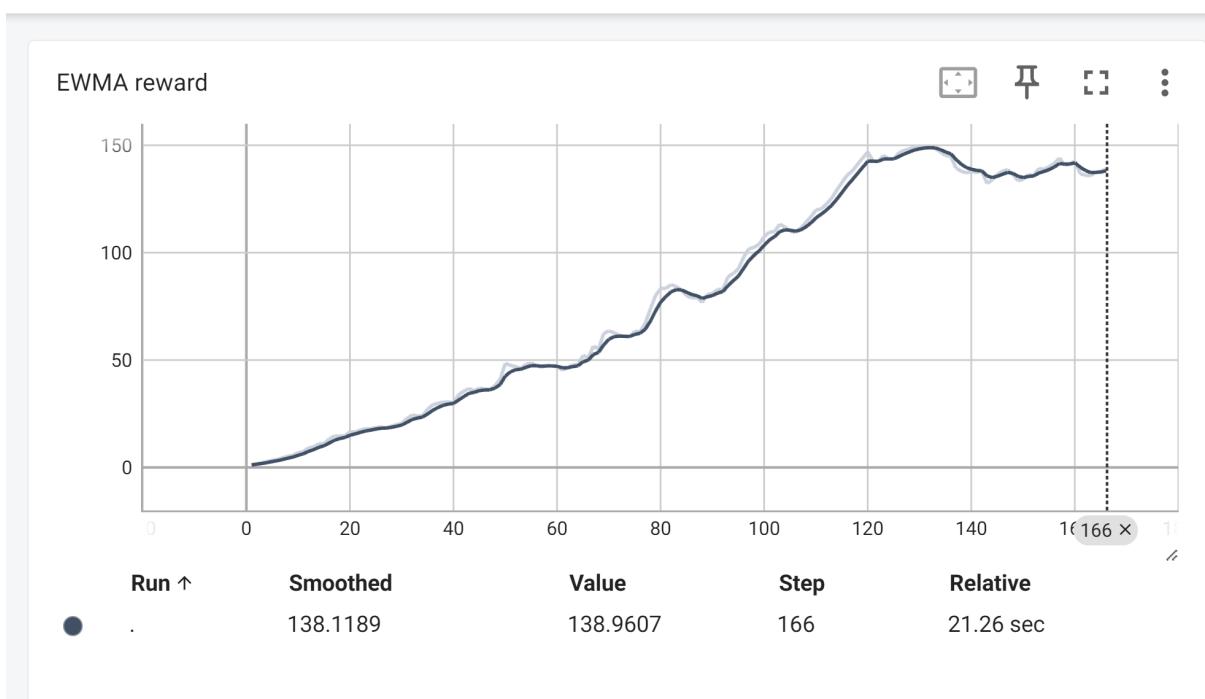
hiddden layer	gamma	learning rate	episode needed
256	0.999	0.01	821

```
▶ Episode 819    length: 200    reward: 200.0    ewma reward: 194.46334559571545
Episode 820    length: 200    reward: 200.0    ewma reward: 194.74017831592968
Episode 821    length: 200    reward: 200.0    ewma reward: 195.0031694001332
Solved! Running reward is now 195.0031694001332 and the last episode runs to 200 time steps!
/usr/local/lib/python3.10/dist-packages/gym/core.py:49: DeprecationWarning: WARN: You are calling render method, but you didn't specified the argument render_mode.
See here for more information: https://www.gymlibrary.ml/content/api/deprecation/
Episode 1    Reward: 200.0
Episode 2    Reward: 200.0
Episode 3    Reward: 200.0
Episode 4    Reward: 200.0
Episode 5    Reward: 200.0
Episode 6    Reward: 200.0
Episode 7    Reward: 200.0
Episode 8    Reward: 200.0
Episode 9    Reward: 200.0
Episode 10   Reward: 200.0
```

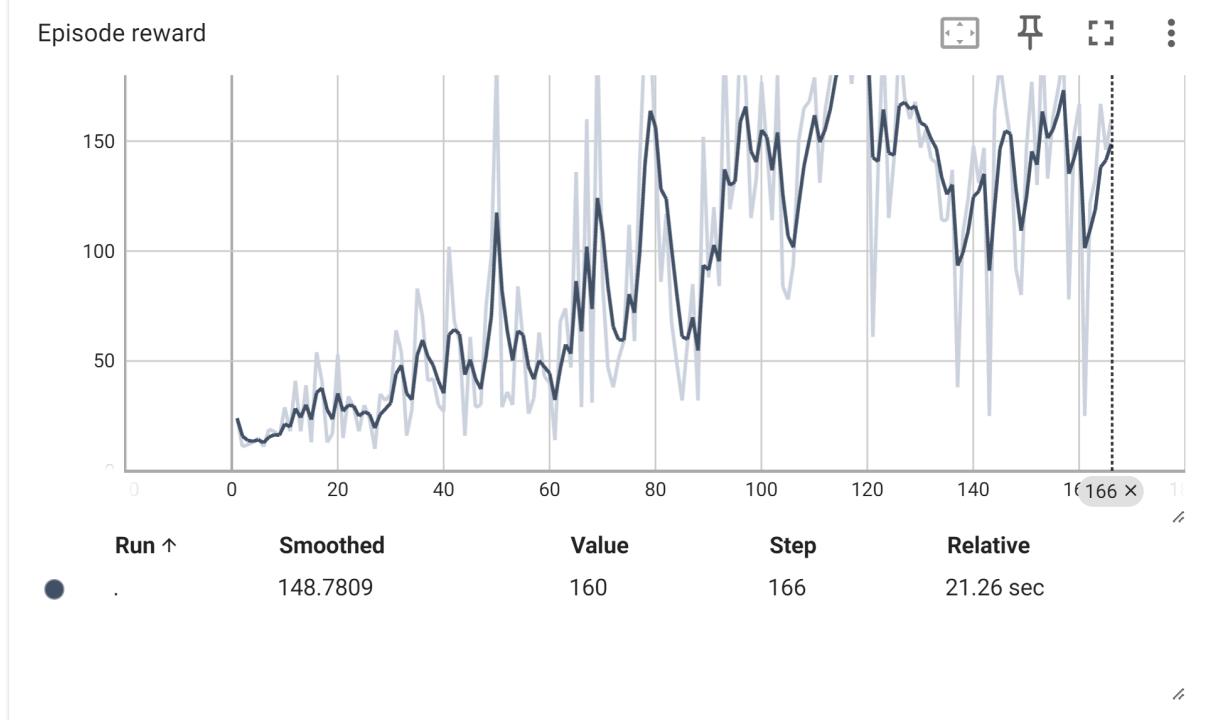


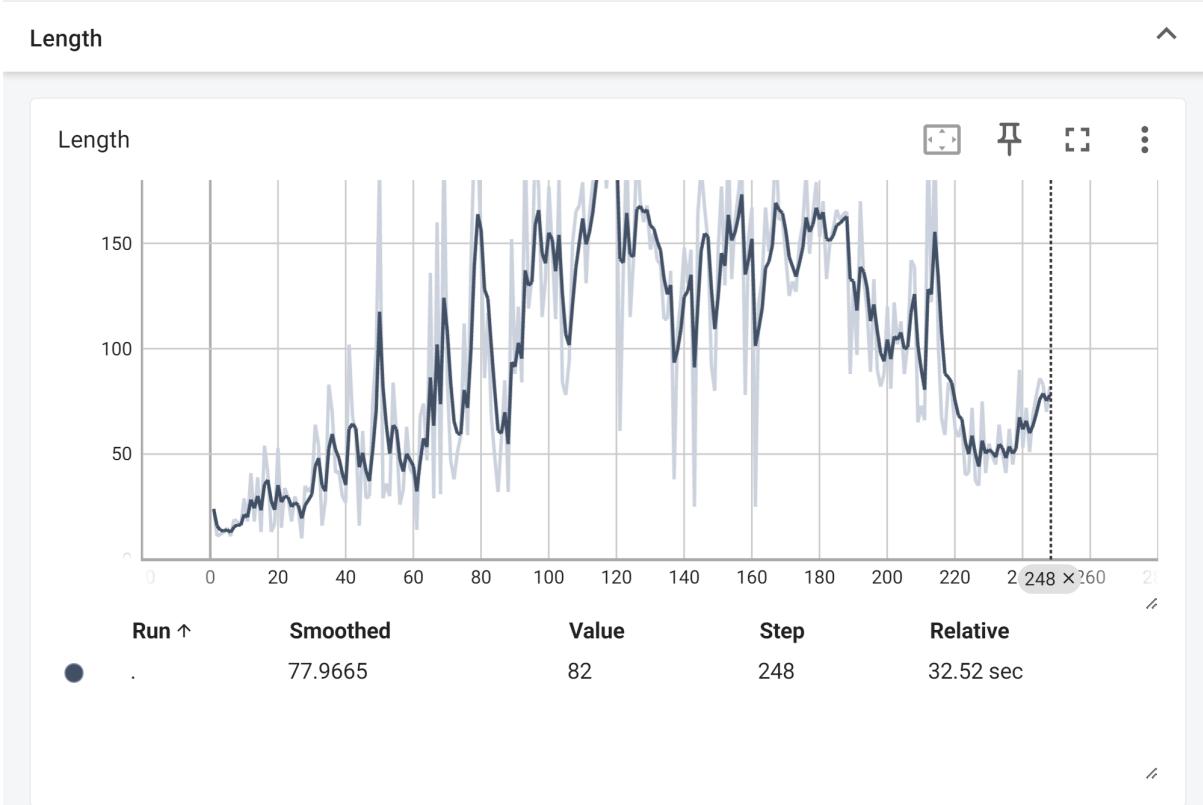
The screenshot shows the TensorBoard interface. At the top, there are tabs for 'TensorBoard', 'TIME SERIES', and 'SCALARS'. The 'SCALARS' tab is active, showing a plot for 'Episode reward'. The x-axis represents episodes from 1 to 166, and the y-axis represents reward values. A single data series is plotted, starting at 148.7809 and ending at 166. The plot area has a light gray background with a white grid. Above the plot, there are two search bars: 'Filter runs (regex)' and 'Filter tags (regex)'. Below the plot, there are several control buttons: 'Run ↑', 'Settings', 'All', 'Scalars', 'Image', 'Histogram', and 'Settings' again. The 'Scalars' button is highlighted. The overall interface is clean with a white background and orange header bar.

EWMA reward



Episode reward





(b)

I add two shared layer in my neuro network, and the activation function are both ReLu. I don't use dropout here.

I found that when the gamma goes down, the episode needed to converge become more. The best number of hidden layer is around 128 - 256. If the learning rate is too large, the model can't converge.

I found that if you don't normalization the return value, you the model can't not converge.

I directly use the Value function as my baseline because I saw that most of the people on the Internet use it, and the performance of Value function as baseline is much better than others baseline created by my brain.

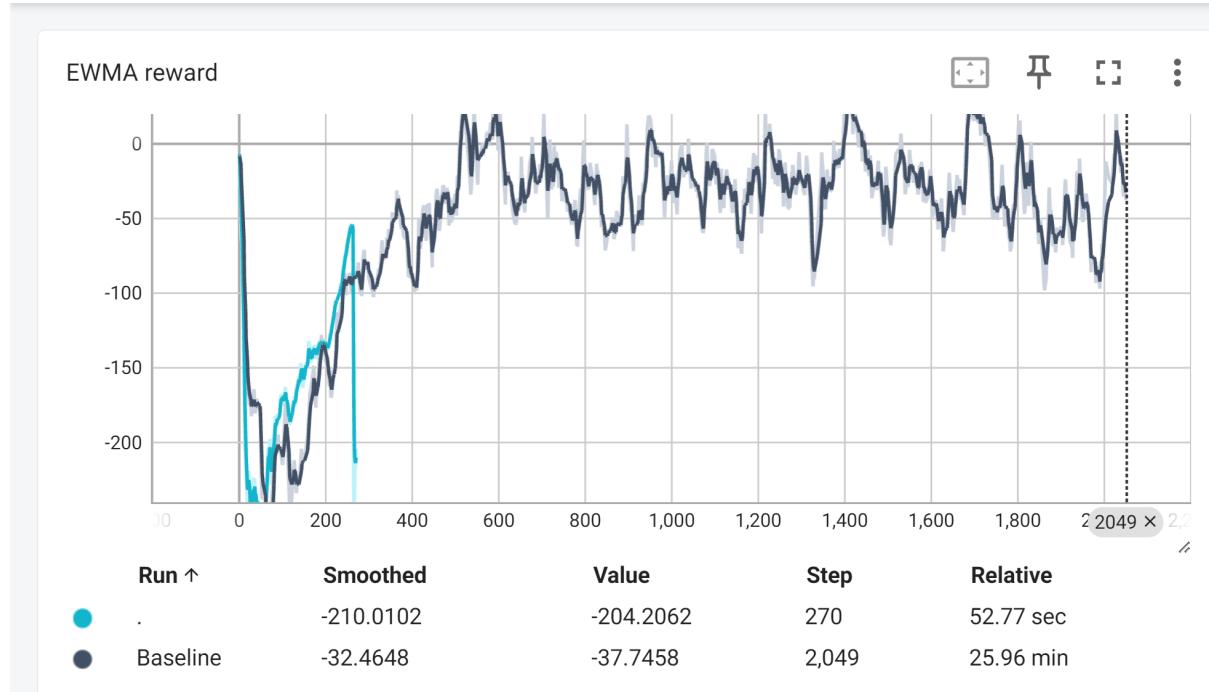
hiddden layer	gamma	learning rate	episode needed
200	0.999	0.02	700

```

Episode 680 length: 191 reward: 42.0/145220000140 ewma reward: 63.30247517700010
Episode 689 length: 291 reward: 41.609508309897734 ewma reward: 64.383824952700776
Episode 690 length: 317 reward: 250.28554953617283 ewma reward: 73.67891118188096
Episode 691 length: 248 reward: 42.695163497003875 ewma reward: 72.1297237976371
Episode 692 length: 590 reward: 129.00051895718104 ewma reward: 74.97326355561428
Episode 693 length: 366 reward: 222.22824313959796 ewma reward: 82.33601253481346
Episode 694 length: 656 reward: 210.3714486727946 ewma reward: 88.73778434171251
Episode 695 length: 405 reward: 253.6177610339569 ewma reward: 96.98178317632473
Episode 696 length: 424 reward: 220.9088162994571 ewma reward: 103.17813483248133
Episode 697 length: 479 reward: 229.13946877145312 ewma reward: 109.47620152942993
Episode 698 length: 417 reward: 153.90794014093436 ewma reward: 111.69778846000516
Episode 699 length: 648 reward: 193.97304454563888 ewma reward: 115.81155126428683
Episode 700 length: 638 reward: 240.39862393505047 ewma reward: 122.040990489782502
Solved! Running reward is now 122.040990489782502 and the last episode runs to 638 time steps!
/usr/local/lib/python3.10/dist-packages/gym/core.py:49: DeprecationWarning: WARN: You are calling render method, but you didn't specified the argument render_m
If you want to render in human mode, initialize the environment in this way: gym.make('EnvName', render_mode='human') and don't call the render method.
See here for more information: https://www.gymlibrary.ml/content/api/
deprecation(
Episode 1 Reward: 196.5985449644424
Episode 2 Reward: -25.48382761693746
Episode 3 Reward: 210.03617741709996
Episode 4 Reward: 240.3397107606958
Episode 5 Reward: 213.14751780437075
Episode 6 Reward: 159.83732163576428
Episode 7 Reward: 232.70809434744083
Episode 8 Reward: 194.098637949563637
Episode 9 Reward: 142.92447458829315
Episode 10 Reward: 231.16624599209732

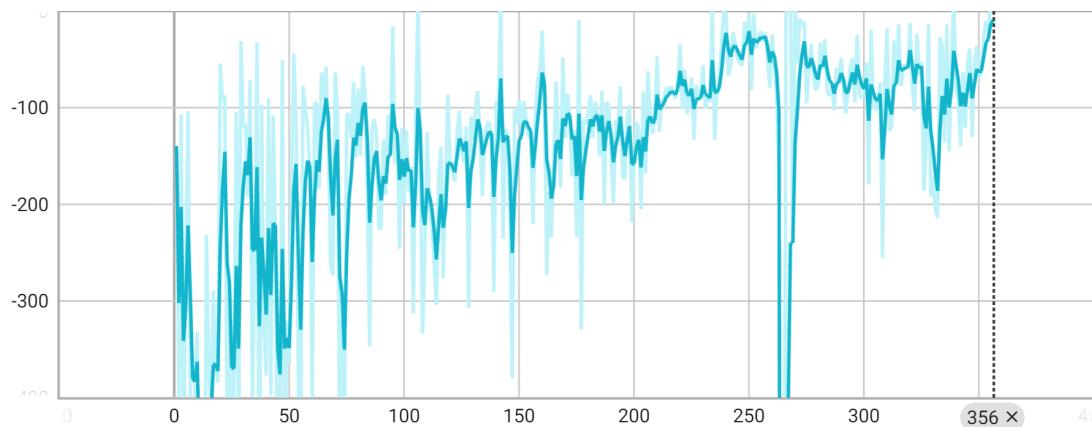
```

EWMA reward



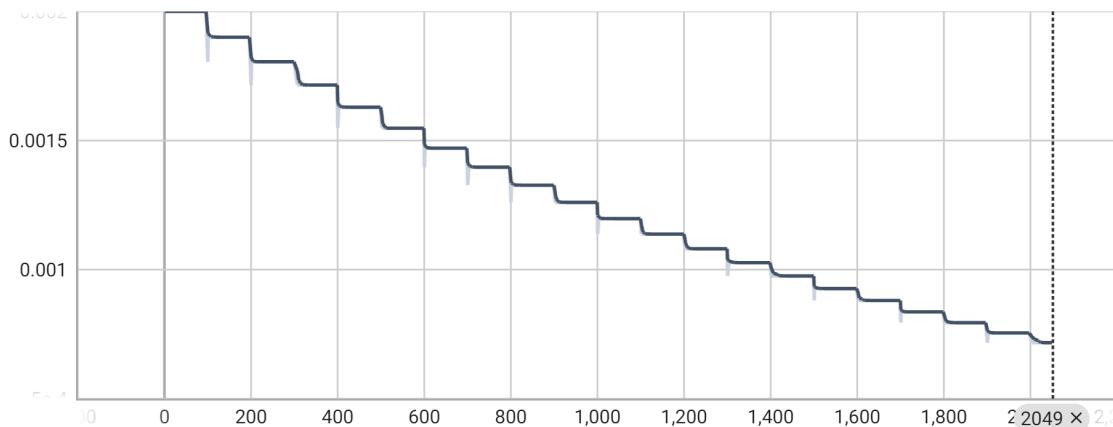
Episode reward

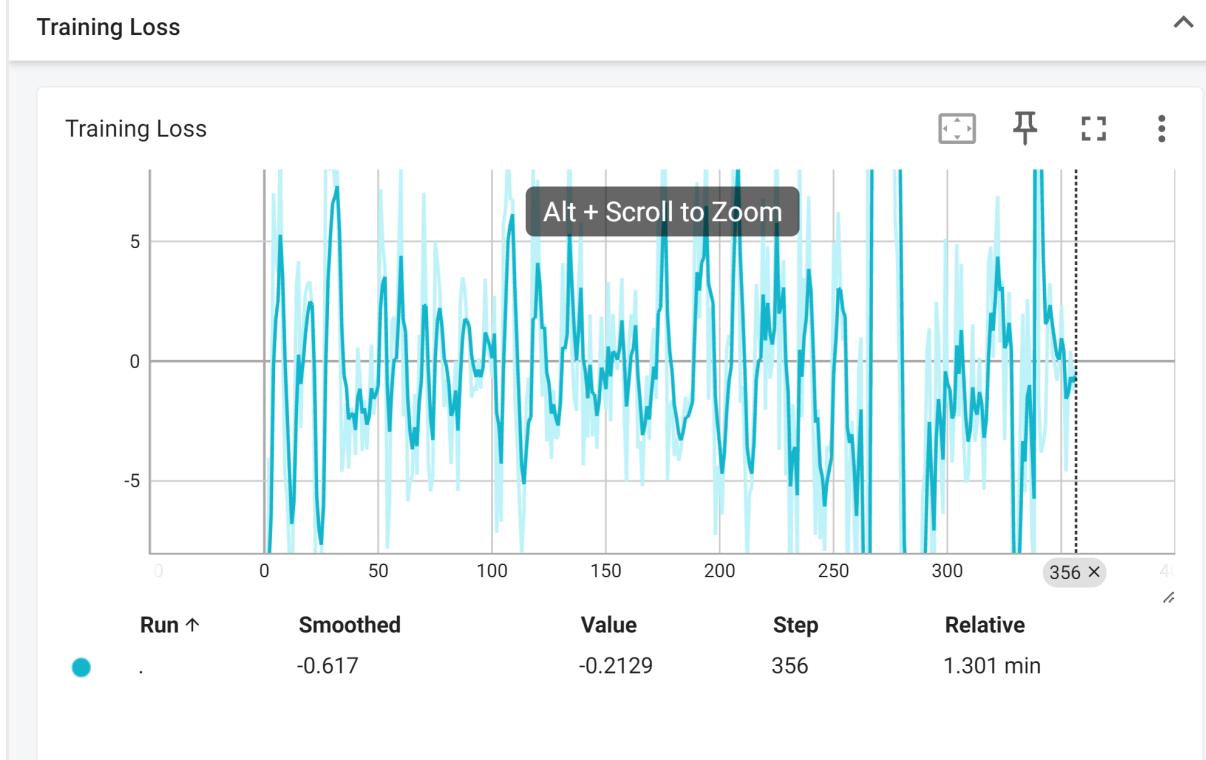
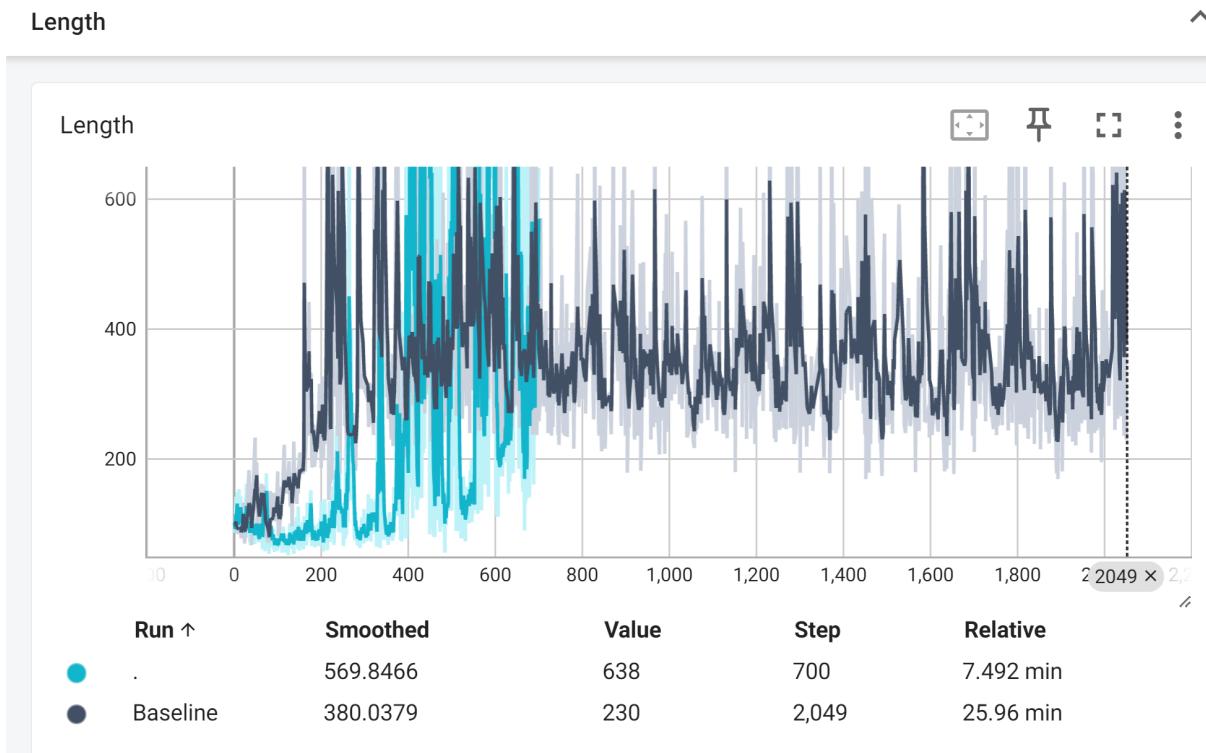
Episode reward



Learning Rate

Learning Rate





(c)

I add two shared layer in my neuro network, and the activation function are both ReLu. I don't use dropout here.

I found that when the gamma goes down, the episode needed to converge become more. The best number of hidden layer is around 128 - 256. If the learning rate is too large, the model can't converge.

I found that if you don't normalization the return value, you the model can't not converge.

I found that if the lambda value is too small, it becomes much more harder to converge. I have tried lambda = 0.5, and although the episode grows to 10000, it hadn't converge. As you can see in the following picure below, the big lambda value require more episode to converge.

lambda = 0.999

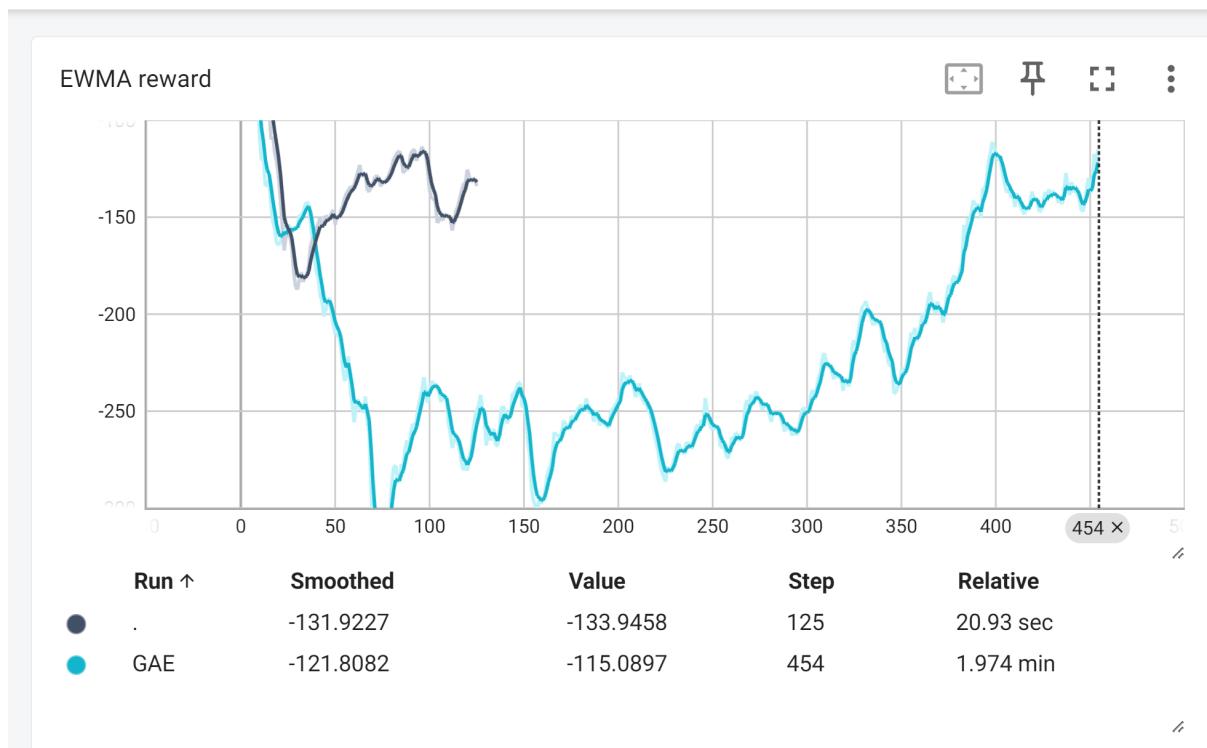
hiddden layer	gamma	learning rate	lambda	episode needed
128	0.999	0.02	0.999	991

```

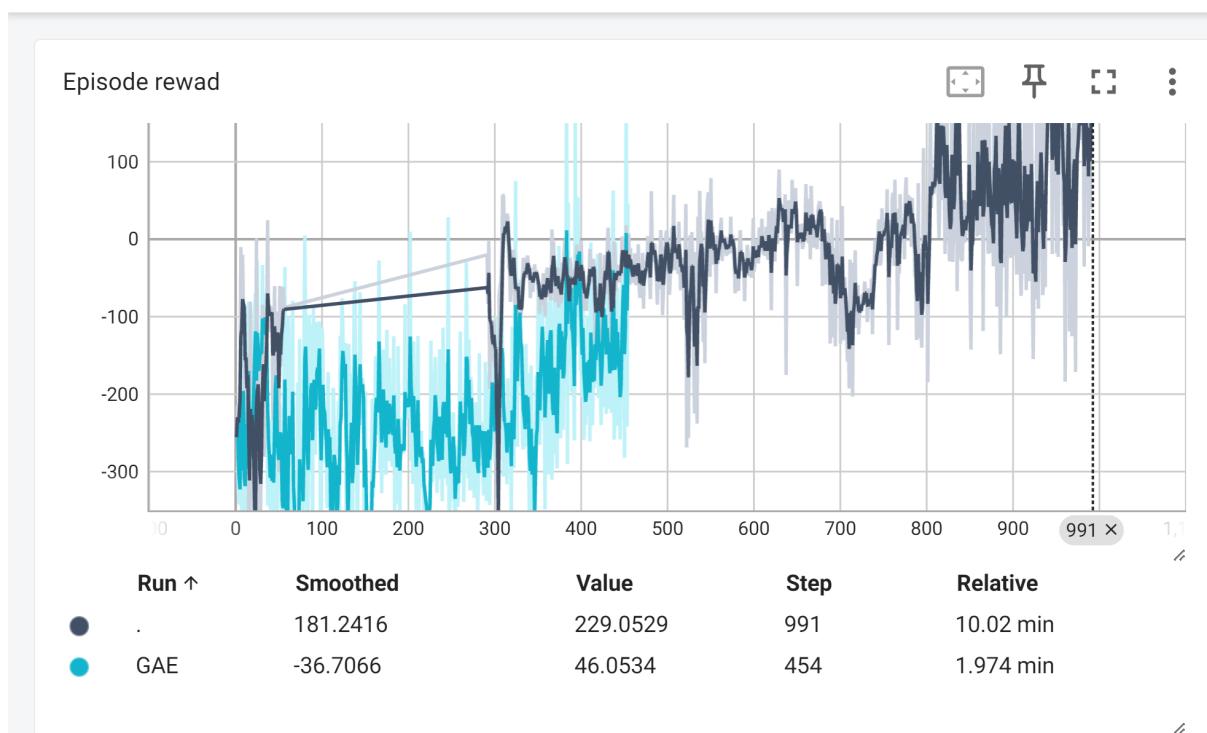
Episode 974 length: 510 reward: 189.02468647594657 ewma reward: 80.69886187841036
Episode 975 length: 735 reward: 144.2105305077076 ewma reward: 83.87444530987521
Episode 976 length: 565 reward: 183.41215531368823 ewma reward: 88.85133081006586
Episode 977 length: 457 reward: 177.92791015005682 ewma reward: 93.30515977706541
Episode 978 length: 636 reward: 128.23865190689764 ewma reward: 95.0518343835702
Episode 979 length: 468 reward: 185.56345186737983 ewma reward: 99.57741525774816
Episode 980 length: 533 reward: 251.25945417964465 ewma reward: 107.16151720384298
Episode 981 length: 758 reward: 154.819534502473 ewma reward: 109.54441806877446
Episode 982 length: 414 reward: -35.35073228711242 ewma reward: 102.29966055098011
Episode 983 length: 605 reward: 155.54759061840082 ewma reward: 104.96205705435113
Episode 984 length: 568 reward: 198.2579583478523 ewma reward: 109.62685211902618
Episode 985 length: 663 reward: 134.37448917789726 ewma reward: 110.86423397196974
Episode 986 length: 769 reward: 137.4608181848917 ewma reward: 112.19406318261582
Episode 987 length: 331 reward: -9.1279184458146 ewma reward: 106.1279641011943
Episode 988 length: 583 reward: 234.4044884500602 ewma reward: 112.54179031863758
Episode 989 length: 266 reward: 38.75072584283291 ewma reward: 108.85223709484734
Episode 990 length: 498 reward: 221.9270852886358 ewma reward: 114.50597950453677
Episode 991 length: 668 reward: 229.05294490445988 ewma reward: 120.23332777453291
Solved! Running reward is now 120.23332777453291 and the last episode runs to 660 time steps!
Episode 1 Reward: 209.33887503677477
Episode 2 Reward: -15.40489345366359
Episode 3 Reward: 162.82639343497598
Episode 4 Reward: 123.37926464268543
Episode 5 Reward: 238.9275000757813
Episode 6 Reward: 165.20897533675662
Episode 7 Reward: -202.39344194732544
Episode 8 Reward: 197.6434686248573
Episode 9 Reward: 153.55577212815058
Episode 10 Reward: 193.14567038029531

```

EWMA reward



Episode reward



lambda = 0.8

hidden layer	gamma	learning rate	lambda	episode needed
--------------	-------	---------------	--------	----------------

128	0.999	0.02	0.8	1121

```

Episode 1109 length: 268 reward: -2.1990687801162494 ewma reward: 78.13826421349869
Episode 1110 length: 319 reward: 257.04478432324817 ewma reward: 87.08359821898616
Episode 1111 length: 683 reward: 276.35733249192134 ewma reward: 96.54727733263292
Episode 1112 length: 572 reward: -57.97371374585377 ewma reward: 88.82122777870858
Episode 1113 length: 170 reward: 30.15199496943107 ewma reward: 85.8877661381203
Episode 1114 length: 210 reward: 5.320178056627967 ewma reward: 81.85938673404569
Episode 1115 length: 546 reward: 226.13028461859477 ewma reward: 89.07293162827314
Episode 1116 length: 349 reward: 245.6056284469832 ewma reward: 96.89956646928864
Episode 1117 length: 272 reward: 43.93170538494422 ewma reward: 94.2511734149954
Episode 1118 length: 333 reward: 291.5794074600534 ewma reward: 104.11758511724831
Episode 1119 length: 313 reward: 281.84165732882417 ewma reward: 113.0037887278271
Episode 1120 length: 397 reward: 214.8232475401364 ewma reward: 118.09476166844256
Episode 1121 length: 451 reward: 233.01837359097956 ewma reward: 123.8409422645694

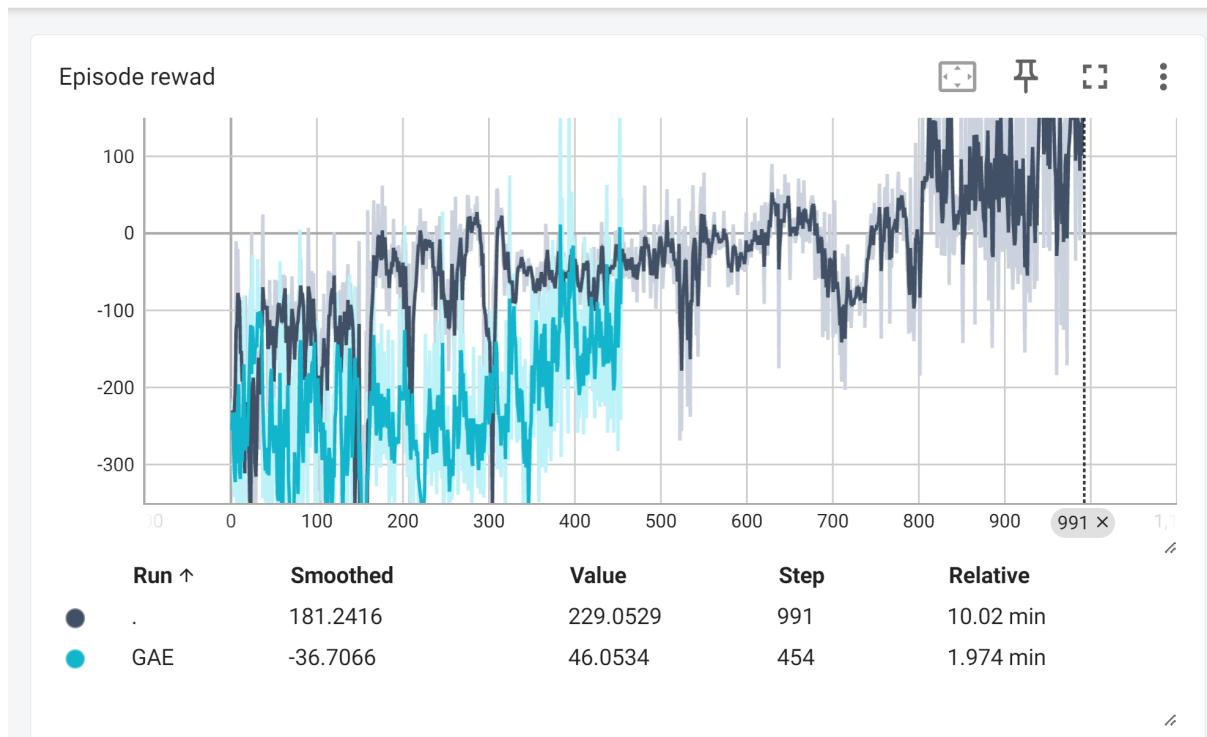
Solved! Running reward is now 123.8409422645694 and the last episode runs to 451 time steps!
/usr/local/lib/python3.10/dist-packages/gym/core.py:49: DeprecationWarning: WARN: You are calling render method, but you didn't specified the argument render_r
If you want to render in human mode, initialize the environment in this way: gym.make('EnvName', render_mode='human') and don't call the render method.
See here for more information: https://www.gymlibrary.ml/content/api/deprecation/
Episode 1 Reward: 247.0316636598362
Episode 2 Reward: 229.62821142472328
Episode 3 Reward: 223.99276251029988
Episode 4 Reward: 25.098684707868827
Episode 5 Reward: -87.95116761930181
Episode 6 Reward: 218.0064286611953
Episode 7 Reward: 281.23469829704686
Episode 8 Reward: 214.73610650776502
Episode 9 Reward: 121.92938109999357
Episode 10 Reward: -5.497648375451604

```

EWMA reward



Episode reward



lambda = 0.7

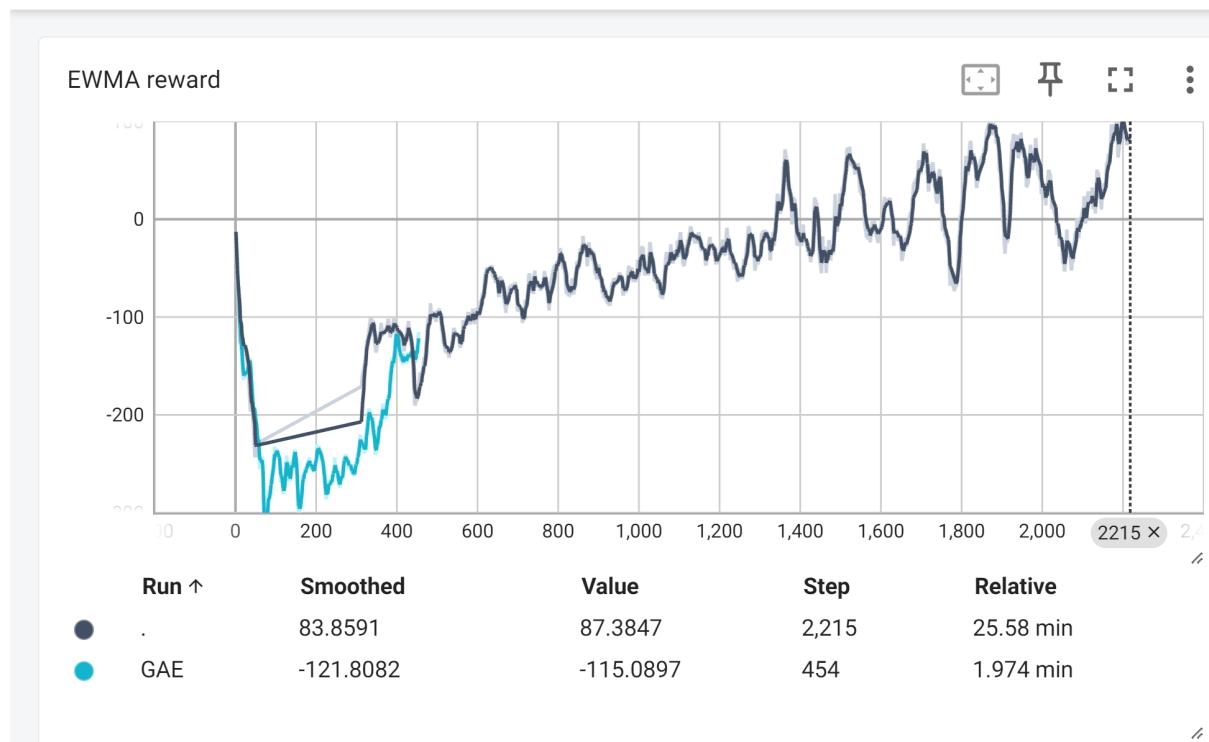
hidden layer	gamma	learning rate	lambda	episode needed
128	0.999	0.02	0.7	2476

```

Episode 2462 length: 310 reward: 290.19495643991803 ewma reward: 49.081531996575386
Episode 2463 length: 266 reward: 274.93390514705203 ewma reward: 60.37415065409921
Episode 2464 length: 221 reward: -153.69216235960477 ewma reward: 49.670835003414005
Episode 2465 length: 363 reward: 22.17562643233373 ewma reward: 48.29607457485999
Episode 2466 length: 403 reward: 234.34283204580527 ewma reward: 57.59841244840725
Episode 2467 length: 397 reward: 246.37512616003065 ewma reward: 67.03724813398841
Episode 2468 length: 369 reward: 247.52191982691144 ewma reward: 76.06148171863455
Episode 2469 length: 344 reward: 245.5563548976075 ewma reward: 84.53622537758321
Episode 2470 length: 418 reward: 228.44987440524326 ewma reward: 91.73190782896621
Episode 2471 length: 388 reward: 238.45261939867464 ewma reward: 99.06794340745164
Episode 2472 length: 338 reward: 270.84159383275266 ewma reward: 107.6566259287167
Episode 2473 length: 387 reward: 231.959651228703 ewma reward: 113.871777193716
Episode 2474 length: 176 reward: 27.291134333496856 ewma reward: 109.54274505070504
Episode 2475 length: 331 reward: 271.27647731765416 ewma reward: 117.62943166405249
Episode 2476 length: 570 reward: 173.3700016465018 ewma reward: 120.41646016317496
Solved! Running reward is now 120.41646016317496 and the last episode runs to 570 time steps!
/usr/local/lib/python3.10/dist-packages/gym/core.py:49: DeprecationWarning: WARN: You are calling render method, but you didn't specified the argument render.
If you want to render in human mode, initialize the environment in this way: gym.make('EnvName', render_mode='human') and don't call the render method.
See here for more information: https://www.gymlibrary.ml/content/api/
deprecation(
Episode 1 Reward: 273.934688881615
Episode 2 Reward: 4.750339050426021
Episode 3 Reward: 241.03648524614556
Episode 4 Reward: 230.78047857005654
Episode 5 Reward: 305.57814213651943
Episode 6 Reward: 199.7426668547442
Episode 7 Reward: 234.50587611265348
Episode 8 Reward: -69.17246394382774
Episode 9 Reward: -2.688902531771703
Episode 10 Reward: 228.81462087743918

```

EWMA reward



Episode reward

