# HW3

## Problem 1

i. **_Proof_**

$$L_{\pi_{\theta_1}}(\pi_{\theta_1}) = \eta(\pi_{\theta_1}) + \sum_{s\in\mathcal{S}} d_\mu^{\pi_{\theta_1}}(s) \sum_{a\in\mathcal{A}} \pi_{\theta_1}(s|a) A^{\pi_{\theta_1}}(s,a)$$

$$= \eta(\pi_{\theta_1}) + \sum_{s\in\mathcal{S}} d_\mu^{\pi_{\theta_1}}(s) \sum_{a\in\mathcal{A}} \pi_{\theta_1}(s|a)(Q^{\pi_{\theta_1}}(s,a) - V^{\pi_{\theta_1}}(s))$$

$$= \eta(\pi_{\theta_1}) + \sum_{s\in\mathcal{S}} d_\mu^{\pi_{\theta_1}}(s) \Big(\sum_{a\in\mathcal{A}} \pi_{\theta_1}(s|a)Q^{\pi_{\theta_1}}(s,a) - V^{\pi_{\theta_1}}(s)\Big)$$

$$= \eta(\pi_{\theta_1}) + \sum_{s\in\mathcal{S}} d_\mu^{\pi_{\theta_1}}(s) \cdot 0$$

ii. To calculate the gradient of $L_{\pi_{\theta_1}}(\pi_\theta)$, differentiate the above expression with respect to $\theta$:

$$\nabla_\theta L_{\pi_{\theta_1}}(\pi_\theta) = \nabla_\theta\left(\eta(\pi_{\theta_1})\right) + \nabla_\theta\left(\sum_{s\in S} d^{\pi_{\theta_1}}(s) \sum_{a\in A} \pi_\theta(a|s) A^{\pi_{\theta_1}}(s,a)\right)$$

Notice that:

- The term $\eta(\pi_{\theta_1})$ does not depend on $\theta$ since it is evaluated at $\theta_1$, so $\nabla_\theta\eta(\pi_{\theta_1}) = 0$.

- The gradient of the second term needs to be evaluated. Here, we differentiate $\pi_\theta(a|s)$ with respect to $\theta$.

Using the policy gradient theorem:

$$\nabla_\theta\left(\sum_{s\in S} d^{\pi_{\theta_1}}(s) \sum_{a\in A} \pi_\theta(a|s) A^{\pi_{\theta_1}}(s,a)\right) = \sum_{s\in S} d^{\pi_{\theta_1}}(s) \sum_{a\in A} \nabla_\theta\pi_\theta(a|s) A^{\pi_{\theta_1}}(s,a)$$

When $\theta = \theta_1$, the expression simplifies because the expectation over the stationary distribution $d^{\pi_{\theta_1}}$ of the advantage function $A^{\pi_{\theta_1}}$ weighted by $\nabla_\theta\pi_\theta(a|s)$ matches the policy gradient formula for the expected return $\eta(\pi_\theta)$:

$$\nabla_\theta\eta(\pi_\theta)\big|\theta = \theta_1 = \sum_{s\in S} d^{\pi_{\theta_1}}(s) \sum_{a\in A} \nabla_\theta\pi_\theta(a|s) A^{\pi_{\theta_1}}(s,a)\big|_{\theta=\theta_1}$$

Therefore:

$$\nabla_\theta L_{\pi_{\theta_1}}(\pi_\theta)\big|\theta = \theta_1 = \nabla_\theta\eta(\pi_\theta)\big|_{\theta=\theta_1}$$

## Problem 2

a. The Lagrangian $\mathcal{L}(\theta, \lambda)$ for the constrained optimization problem is given by:

$$\mathcal{L}(\theta,\lambda) = -(\nabla_\theta L_{\theta_k}(\theta)|_{\theta=\theta_k})^T(\theta - \theta_k) + \lambda\left(\frac{1}{2}(\theta-\theta_k)^T H_{\theta_k}(\theta - \theta_k) - \delta\right)$$

Setting $\frac{\partial\mathcal{L}}{\partial\theta} = 0$ leads to:

$$-\nabla_\theta L_{\theta_k}(\theta)|_{\theta=\theta_k} + \lambda H_{\theta_k}(\theta - \theta_k) = 0$$

Then we can get $\theta$:

$$\theta = \theta_k + \frac{1}{\lambda}H_{\theta_k}^{-1}\nabla_\theta L_{\theta_k}(\theta)|_{\theta=\theta_k}$$

Plug $\theta$ back into $\mathcal{L}(\theta,\lambda)$:

$$\mathcal{L}(\theta,\lambda) = -(\nabla_\theta L_{\theta_k}(\theta)|_{\theta=\theta_k})^T(\frac{1}{\lambda}H_{\theta_k}^{-1}\nabla_\theta L_{\theta_k}(\theta)|_{\theta=\theta_k}) + \lambda\left(\frac{1}{2}(\frac{1}{\lambda}H_{\theta_k}^{-1}\nabla_\theta L_{\theta_k}(\theta)|_{\theta=\theta_k})^T H_{\theta_k}(\frac{1}{\lambda}H_{\theta_k}^{-1}\nabla_\theta L_{\theta_k}(\theta)|_{\theta=\theta_k}) - \delta\right)$$

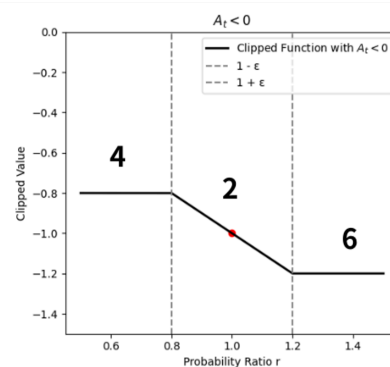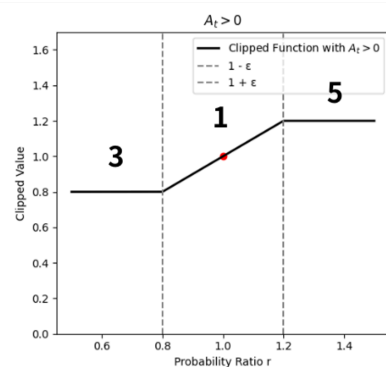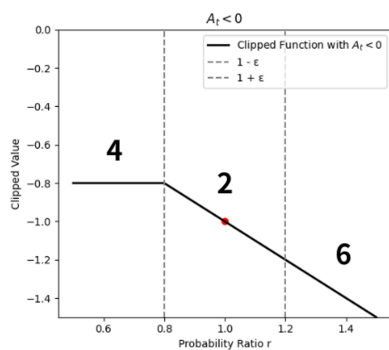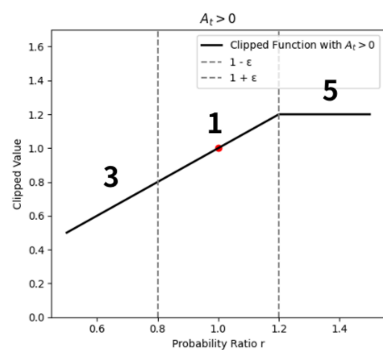$$= -\frac{1}{\lambda}\left(\nabla_\theta L_{\theta_k}(\theta)|_{\theta=\theta_k}\right)^T H_{\theta_k}^{-1}\left(\nabla_\theta L_{\theta_k}(\theta)|_{\theta=\theta_k}\right) + \frac{1}{2\lambda}\left(\nabla_\theta L_{\theta_k}(\theta)|_{\theta=\theta_k}\right)^T H_{\theta_k}^{-1}\left(\nabla_\theta L_{\theta_k}(\theta)|_{\theta=\theta_k}\right) - \lambda\delta$$

$$= -\frac{1}{2\lambda}\left(\nabla_\theta L_{\theta_k}(\theta)|_{\theta=\theta_k}\right)^T H_{\theta_k}^{-1}\left(\nabla_\theta L_{\theta_k}(\theta)|_{\theta=\theta_k}\right) - \lambda\delta$$

Hence, $D(\lambda)$ is as provided. As for $\lambda^*$:

$$\frac{d}{d\lambda}D(\lambda) = \frac{1}{2\lambda^2}\left(\nabla_\theta L_{\theta_k}(\theta)|_{\theta=\theta_k}\right)^T H_{\theta_k}^{-1}\left(\nabla_\theta L_{\theta_k}(\theta)|_{\theta=\theta_k}\right) - \delta = 0$$

$$\lambda^* = \sqrt{\frac{\left(\nabla\theta L_{\theta_k}(\theta)|_{\theta=\theta_k}\right)^T H_{\theta_k}^{-1}\left(\nabla_\theta L_{\theta_k}(\theta)|_{\theta=\theta_k}\right)}{2\delta}}$$

b. In part a, we already prove that $\theta^* = \theta_k + \frac{1}{\lambda}H_{\theta_k}^{-1}\nabla_\theta L_{\theta_k}(\theta)|_{\theta=\theta_k}$ if $\theta^* = \arg\min_\theta \mathcal{L}(\theta, \lambda)$, then we can get $\alpha$:

$$\alpha = \frac{1}{\lambda^*} = \sqrt{\frac{2\delta}{\left(\nabla\theta L_{\theta_k}(\theta)|_{\theta=\theta_k}\right)^T H_{\theta_k}^{-1}\left(\nabla_\theta L_{\theta_k}(\theta)|_{\theta=\theta_k}\right)}}$$

# Problem 3

**Original PPO-Clip**

$$L_{s,a}^{\text{clip}}(\theta;\theta_k) := \min\left\{\frac{\pi_\theta(a\mid s)}{\pi_{\theta_k}(a\mid s)}A^{\theta_k}(s,a), \text{clip}\left(\frac{\pi_\theta(a\mid s)}{\pi_{\theta_k}(a\mid s)}, 1-\epsilon, 1+\epsilon\right)A^{\theta_k}(s,a)\right\}$$

**Variant PPO-Clip**

$$\tilde{L}_{s,a}^{\text{clip}}(\theta;\theta_k) = \text{clip}\left(\frac{\pi_\theta(a\mid s)}{\pi_{\theta_k}(a\mid s)}, 1-\epsilon, 1+\epsilon\right)A^{\theta_k}(s,a)$$

| $p_t(\theta) > 0$ | $A_t$ | Return Value of *min* | Objective is Clipped | Sign of Objective | Gradient |
|---|---|---|---|---|---|
| $p_t(\theta) \in [1-\epsilon, 1+\epsilon]$ | $+$ | $p_t(\theta)A_t$ | no | $+$ | ✓ |
| $p_t(\theta) \in [1-\epsilon, 1+\epsilon]$ | $-$ | $p_t(\theta)A_t$ | no | $-$ | ✓ |
| $p_t(\theta) < 1-\epsilon$ | $+$ | $p_t(\theta)A_t$ | no | $+$ | ✓ |
| $p_t(\theta) < 1-\epsilon$ | $-$ | $(1-\epsilon)A_t$ | yes | $-$ | 0 |
| $p_t(\theta) > 1+\epsilon$ | $+$ | $(1+\epsilon)A_t$ | yes | $+$ | 0 |
| $p_t(\theta) > 1+\epsilon$ | $-$ | $p_t(\theta)A_t$ | no | $-$ | ✓ |

| $p_t(\theta) > 0$ | $A_t$ | Return Value of *min* | Objective is Clipped | Sign of Objective |
|---|---|---|---|---|
| $p_t(\theta) \in [1-\epsilon, 1+\epsilon]$ | $+$ | $p_t(\theta)A_t$ | no | $+$ |
| $p_t(\theta) \in [1-\epsilon, 1+\epsilon]$ | $-$ | $p_t(\theta)A_t$ | no | $-$ |
| $p_t(\theta) < 1-\epsilon$ | $+$ | $(1-\epsilon)A_t$ | yes | $+$ |
| $p_t(\theta) < 1-\epsilon$ | $-$ | $(1-\epsilon)A_t$ | yes | $-$ |
| $p_t(\theta) > 1+\epsilon$ | $+$ | $(1+\epsilon)A_t$ | yes | $+$ |
| $p_t(\theta) > 1+\epsilon$ | $-$ | $(1+\epsilon)A_t$ | yes | $-$ |





The purpose of applying the clipping mechanism is to prevent the policy from updating too significantly in the correct direction. Therefore, there is no need to clip when the updates move in the wrong direction. Taking the minimum of the original value ensures this; thus, the function will continuously decrease.

# Problem 4

## a. Pendulum-v1

- **NN architecture**
  - **Actor**

```
nn.Sequential(
        nn.Linear(num_inputs, 40),
        nn.ReLU(),
        nn.Linear(400, 300),
        nn.ReLU(),
        nn.Linear(300, num_outputs),
        nn.Tanh()
)
```

  - **Critic**

```
nn.Sequential(
    nn.Linear(num_inputs, 400),
    nn.ReLU()
)
nn.Sequential(
    nn.Linear(400 + num_outputs, 300),
    nn.ReLU(),
    nn.Linear(300, 1)
)
```
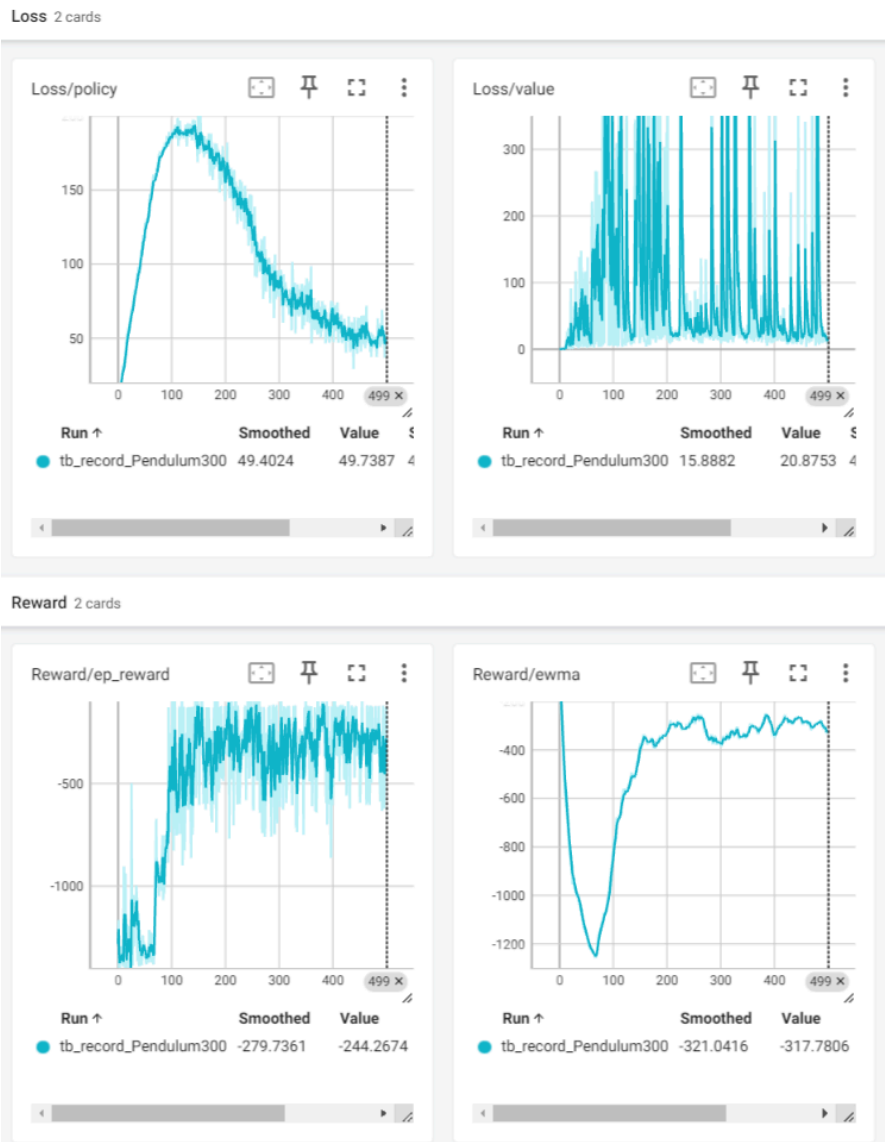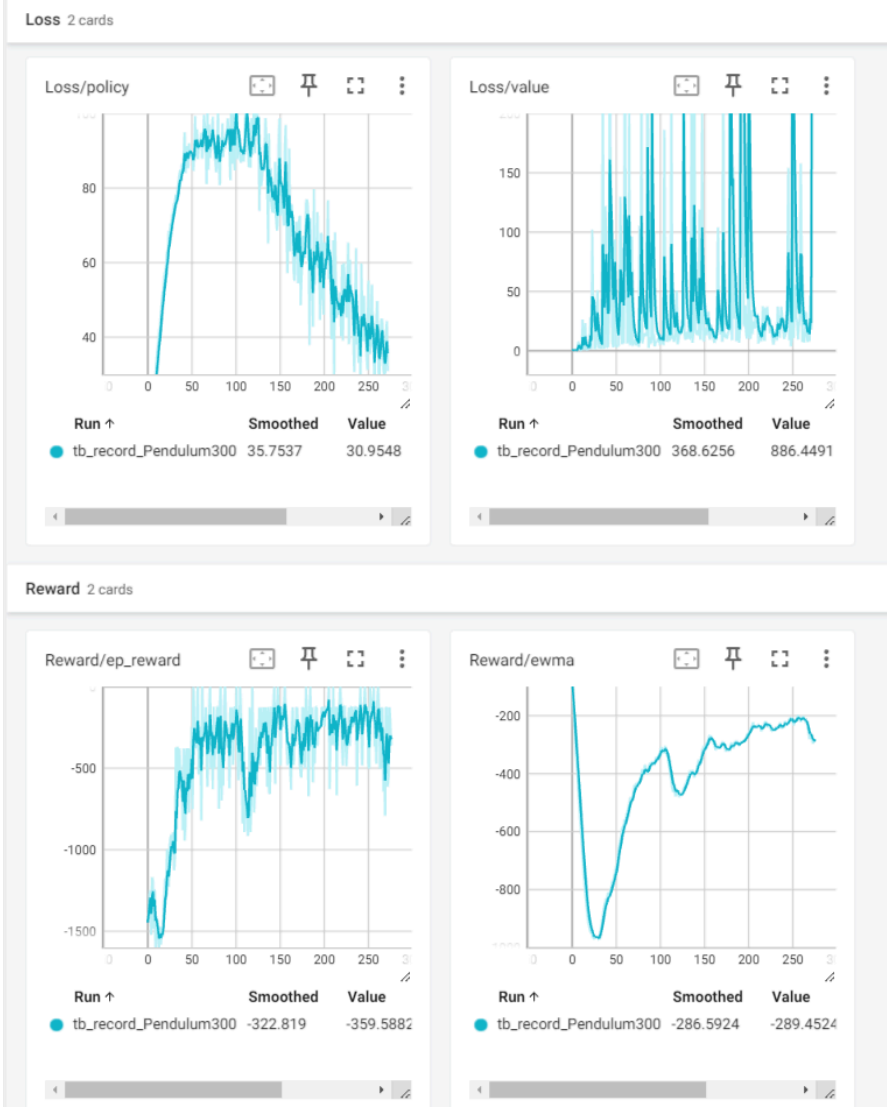
- **Training Result**

  Hyperparameter                                    Hyperparameter

```
num_episodes = 300
gamma = 0.995
tau = 0.002
hidden_size = 128
noise_scale = 0.1
replay_size = 100000
batch_size = 128
updates_per_step = 1
```

```
num_episodes = 500
gamma = 0.995
tau = 0.002
hidden_size = 128
noise_scale = 0.3
replay_size = 100000
batch_size = 128
updates_per_step = 1
```
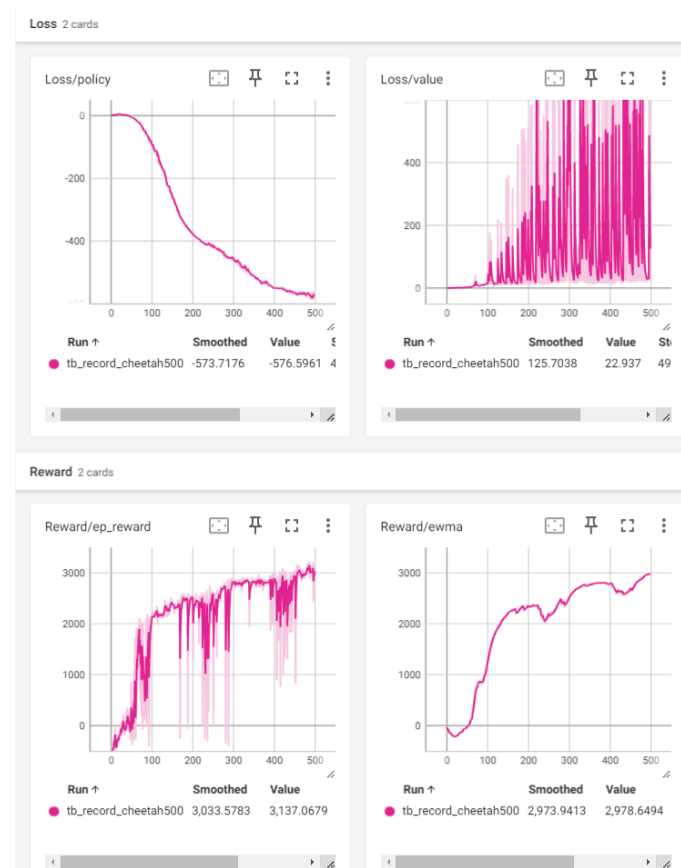


- **Discussion**

  1. **Loss/Policy**:

     This graph shows the loss associated with the actor network over the course of training episodes. The loss peaks around episode 100 and then generally trends downward, indicating that the actor network is gradually learning a strategy that minimizes error in its predictions of actions.

  2. **Loss/Value**:

     This chart shows the loss of the critic network. The frequent spikes in loss indicate periods where the network's predictions are significantly off from actual returns, but a general downward trend or stabilization would be expected as training progresses. However, the volatility here suggests that the value estimates might still be fluctuating significantly, which could be due to several factors like exploration strategies or parameter settings.

  3. **Reward/ep_reward**:

     The reward per episode chart shows the reward obtained by the agent each episode. It's clear from the chart that the reward dramatically increases at the beginning of training, indicating that the agent quickly learns a reasonably effective strategy. The reward stabilizes after an initial learning phase, with some fluctuations which could be due to exploration or environmental stochasticity.

  4. **Reward/ewma**:

The exponentially weighted moving average (EWMA) of the reward smooths out the episodic rewards to give a clearer picture of the overall trend. This curve, similarly to the episodic reward, shows significant improvement in the initial episodes and then levels off, indicating that the policy has likely converged to a near-optimal strategy for this environment. Comparing the two charts above, we can observe that even after an additional 200 training episodes, the EWMA (Exponential Weighted Moving Average) has not shown a significant increase. Therefore, it can be understood that an EWMA of approximately -200 might be the limit value for this environment.

- Encountered issue: Initially, when running `ddpg.py` on a laptop, no matter how the code was modified, the policy could not converge. After transferring the same code to a PC, convergence was achieved. It is speculated that the issue was due to insufficient memory on the laptop, leading to a full replay buffer that could not learn anything new.

## b. Pendulum-v3

- **NN architecture**
  - **Actor**

```
nn.Sequential(
        nn.Linear(num_inputs, 40),
        nn.ReLU(),
        nn.Linear(400, 300),
        nn.ReLU(),
        nn.Linear(300, num_outputs),
        nn.Tanh()
)
```

  - **Critic**

```
nn.Sequential(
    nn.Linear(num_inputs, 400),
    nn.ReLU()
)
nn.Sequential(
    nn.Linear(400 + num_outputs, 300),
    nn.ReLU(),
    nn.Linear(300, 1)
)
```

- **Hyperparameter**

```
gamma = 0.995
tau = 0.002
hidden_size = 128
noise_scale = 0.1
replay_size = 100000
batch_size = 128
updates_per_step = 1
```
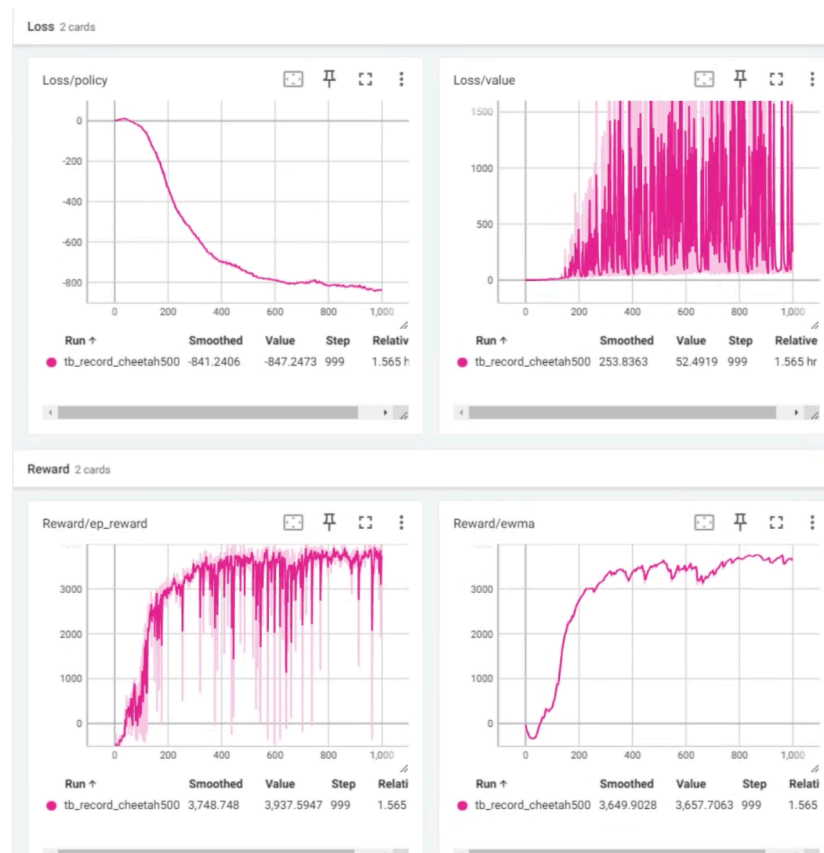
- **Training Result**

Training for 1000 episodes



Training for 500 episodes

- **Discussion**

  The policy loss decreasing significantly over time. The smooth line indicates a consistent reduction in loss, suggesting improvements in policy behavior throughout the training. EWMA of rewards smoothens the fluctuations seen in the per-episode rewards chart. This curve rises steadily, indicating consistent improvement in the model's performance over time, with fewer extremes compared to the per-episode reward graph. We can also observe that before each increase to a new peak in rewards, the EWMA of the rewards first dips. This could be due to the policy escaping from a local maximum.