

Online and Reinforcement Learning (2025)

Home Assignment 6

Davide Marchi 777881

Contents

1	PPO	2
1.1	Return expressed as advantage over another policy	2
1.2	Clipping	3
1.3	Pi prime in PPO	4
2	Offline Evaluation of Bandit Algorithms	4
2.1	Part 1	4
2.2	Part 2	5

1 PPO

1.1 Return expressed as advantage over another policy

We wish to prove that the expected return of a policy π can be written as

$$J(\pi) = J(\pi_{\text{ref}}) + \mathbb{E}_{s_0 \sim p_0, \pi} \left[\sum_{t=0}^{\infty} \gamma^t A^{\pi_{\text{ref}}}(s_t, a_t) \right], \quad (1)$$

where the advantage function is defined by

$$A^{\pi_{\text{ref}}}(s, a) = Q^{\pi_{\text{ref}}}(s, a) - V^{\pi_{\text{ref}}}(s).$$

Proof: Recall that for any state s the state-value function of π_{ref} is given by

$$V^{\pi_{\text{ref}}}(s) = \mathbb{E}_{a \sim \pi_{\text{ref}}} [Q^{\pi_{\text{ref}}}(s, a)].$$

Consider a trajectory $\tau = (s_0, a_0, s_1, a_1, \dots)$ generated by following π . For any finite horizon T , we can write

$$\begin{aligned} \sum_{t=0}^T \gamma^t A^{\pi_{\text{ref}}}(s_t, a_t) &= \sum_{t=0}^T \gamma^t (Q^{\pi_{\text{ref}}}(s_t, a_t) - V^{\pi_{\text{ref}}}(s_t)) \\ &= \sum_{t=0}^T \gamma^t (r(s_t, a_t) + \gamma V^{\pi_{\text{ref}}}(s_{t+1}) - V^{\pi_{\text{ref}}}(s_t)). \end{aligned}$$

Notice that the sum telescopes. To see this, rearrange the terms:

$$\sum_{t=0}^T \gamma^t r(s_t, a_t) + \sum_{t=0}^T (\gamma^{t+1} V^{\pi_{\text{ref}}}(s_{t+1}) - \gamma^t V^{\pi_{\text{ref}}}(s_t)).$$

The second sum is telescopic:

$$\sum_{t=0}^T (\gamma^{t+1} V^{\pi_{\text{ref}}}(s_{t+1}) - \gamma^t V^{\pi_{\text{ref}}}(s_t)) = -V^{\pi_{\text{ref}}}(s_0) + \gamma^{T+1} V^{\pi_{\text{ref}}}(s_{T+1}).$$

Assuming that $V^{\pi_{\text{ref}}}(s)$ is bounded and $\gamma \in (0, 1)$, as $T \rightarrow \infty$ we have $\gamma^{T+1} V^{\pi_{\text{ref}}}(s_{T+1}) \rightarrow 0$. Therefore,

$$\sum_{t=0}^{\infty} \gamma^t A^{\pi_{\text{ref}}}(s_t, a_t) = \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) - V^{\pi_{\text{ref}}}(s_0).$$

Taking the expectation over trajectories starting from $s_0 \sim p_0$ (and using the definition $J(\pi) = \mathbb{E}[\sum_{t \geq 0} \gamma^t r(s_t, a_t)]$) gives

$$J(\pi) = \mathbb{E}_{s_0 \sim p_0, \pi} \left[\sum_{t=0}^{\infty} \gamma^t A^{\pi_{\text{ref}}}(s_t, a_t) \right] + \mathbb{E}_{s_0 \sim p_0} [V^{\pi_{\text{ref}}}(s_0)].$$

Since by definition $J(\pi_{\text{ref}}) = \mathbb{E}_{s_0 \sim p_0} [V^{\pi_{\text{ref}}}(s_0)]$, we obtain the desired result (1).

1.2 Clipping

In PPO the surrogate objective for a given state-action pair is defined as

$$L^{\text{CLIP}}(\theta) = \min \left(r(\theta) \hat{A}^{\pi_{\text{ref}}}(s, a), \text{clip}(r(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}^{\pi_{\text{ref}}}(s, a) \right),$$

where

$$r(\theta) = \frac{\pi_{\theta}(a|s)}{\pi_{\text{ref}}(a|s)}$$

and the clipping function is defined by

$$\text{clip}(x, l, u) = \min\{\max(x, l), u\}.$$

Intuitive discussion: The gradient update will change $\pi_{\theta}(a|s)$ (and hence $r(\theta)$) if either

- (i) $r(\theta) \in [1 - \epsilon, 1 + \epsilon]$, or
- (ii) when $r(\theta) \notin [1 - \epsilon, 1 + \epsilon]$, the *unclipped* term $r(\theta) \hat{A}^{\pi_{\text{ref}}}(s, a)$ has a gradient that *points toward* the interval $[1 - \epsilon, 1 + \epsilon]$.

The first condition is trivial. Now assume that $r(\theta) \notin [1 - \epsilon, 1 + \epsilon]$. There are two cases:

Case 1: $r(\theta) > 1 + \epsilon$. Then

$$\text{clip}(r(\theta), 1 - \epsilon, 1 + \epsilon) = 1 + \epsilon.$$

Now, consider the sign of $\hat{A}^{\pi_{\text{ref}}}(s, a)$:

- If $\hat{A}^{\pi_{\text{ref}}}(s, a) > 0$, then

$$r(\theta) \hat{A}^{\pi_{\text{ref}}}(s, a) > (1 + \epsilon) \hat{A}^{\pi_{\text{ref}}}(s, a).$$

Thus, the minimum in the objective is the clipped term, which is constant with respect to θ (i.e., its gradient is zero).

- If $\hat{A}^{\pi_{\text{ref}}}(s, a) < 0$, then

$$r(\theta) \hat{A}^{\pi_{\text{ref}}}(s, a) < (1 + \epsilon) \hat{A}^{\pi_{\text{ref}}}(s, a),$$

so the unclipped term is active. Its gradient with respect to θ is proportional to

$$\nabla_{\theta} \left[r(\theta) \hat{A}^{\pi_{\text{ref}}}(s, a) \right] = \hat{A}^{\pi_{\text{ref}}}(s, a) \nabla_{\theta} r(\theta).$$

Since $\hat{A}^{\pi_{\text{ref}}}(s, a) < 0$, the gradient will be *negative* (assuming $\nabla_{\theta} r(\theta) > 0$ for an increase in $r(\theta)$), which implies that the update will *decrease* $r(\theta)$ — that is, it pushes $r(\theta)$ *toward* the boundary $1 + \epsilon$ rather than away from it.

Case 2: $r(\theta) < 1 - \epsilon$. By a similar argument, if $\hat{A}^{\pi_{\text{ref}}}(s, a) > 0$, then the gradient of the unclipped term (which is active in this case) is positive and pushes $r(\theta)$ upward toward $1 - \epsilon$.

Formalization: Assume that $r(\theta) \notin [1 - \epsilon, 1 + \epsilon]$. Then:

$$\begin{cases} r(\theta) > 1 + \epsilon & \text{and} & \hat{A}^{\pi_{\text{ref}}}(s, a) < 0, \\ r(\theta) < 1 - \epsilon & \text{and} & \hat{A}^{\pi_{\text{ref}}}(s, a) > 0. \end{cases}$$

In these cases, the derivative of the unclipped objective is

$$\nabla_{\theta} \left[r(\theta) \hat{A}^{\pi_{\text{ref}}}(s, a) \right] = \hat{A}^{\pi_{\text{ref}}}(s, a) \nabla_{\theta} r(\theta).$$

Thus, when $\hat{A}^{\pi_{\text{ref}}}(s, a)$ has the sign that makes this derivative point toward the interval, the gradient-based update will reduce the deviation of $r(\theta)$ from $[1 - \epsilon, 1 + \epsilon]$. In other words, even if the current ratio is outside the interval, the gradient direction (if it does not point away from the interval) will drive it closer to the interval, ensuring that the policy update is conservative.

1.3 Pi prime in PPO

In the *Gather experience* phase, the policy that generates the data is π with parameters θ' (denoted as the behavior policy), and the probability of taking action a_t^e in state s_t^e is stored as

$$p_t^e = \pi_{\theta'}(a_t^e | s_t^e).$$

Later, during the PPO update the current policy π_{θ} (with updated parameters θ) is used. Thus, the ratio

$$\frac{\pi_{\theta}(a_t^e | s_t^e)}{p_t^e} = \frac{\pi_{\theta}(a_t^e | s_t^e)}{\pi_{\theta'}(a_t^e | s_t^e)}$$

generally differs from 1 because $\theta \neq \theta'$ in general. In other words, since the policy is updated over time, the probability of taking the same action in the same state under the current policy is typically not equal to the probability under the behavior policy that generated the data. This ratio is used as an importance sampling correction to account for the discrepancy between the two policies.

2 Offline Evaluation of Bandit Algorithms

2.1 Part 1

Evaluating algorithms for online learning with limited feedback in real-life scenarios is challenging. While the most direct approach is to implement an algorithm and measure its performance live, this is often not feasible due to:

- **Potential risks, costs, and time delays.** Deploying a potentially suboptimal algorithm can lead to high financial or reputational costs. Moreover, once the algorithm is running, it takes time to collect sufficient data for analysis, which delays insights and can be inefficient if the algorithm underperforms.
- **Difficulty of controlled experimentation.** Once data is collected based on a particular algorithm’s actions, it is nearly impossible to “replay” the same conditions to test different algorithms under identical circumstances. This lack of controlled repetition makes fair comparisons difficult.

2.2 Part 2

(a) Modification of UCB1 for Importance-Weighted Losses (Uniform Sampling)

To handle partial feedback when arms are chosen uniformly at random (with probability $1/K$ for each arm), we replace the usual empirical loss estimates in UCB1 with importance-weighted estimates. Concretely, whenever an arm i is chosen, its observed loss is scaled by the factor $\frac{1}{p_i(t)} = K$, ensuring an unbiased estimator. The main changes from standard UCB1 are thus:

- **Empirical Loss Update:** Instead of adding the raw observed loss $\ell_{i,t}$, we add $K \cdot \ell_{i,t}$ to the running total for arm i .
- **Confidence Bounds:** The increased variance (due to multiplying by K) is accounted for in the confidence term, typically by a constant factor in front of the usual $\sqrt{\frac{\ln t}{N_i(t)}}$ bound.

Pseudo-Code of the Modified Algorithm:

Initialize: For each arm $i \in \{1, \dots, K\}$,

$$\widehat{L}_i(0) = 0, \quad N_i(0) = 0.$$

For $t = 1$ **to** T :

1. Select arm

$$A_t = \arg \min_{i \in \{1, \dots, K\}} \left(\widehat{L}_i(t-1) + c \sqrt{\frac{\ln(t-1)}{N_i(t-1)}} \right),$$

where $c > 0$ is a constant.

2. Observe the loss $\ell_{A_t,t}$ for the chosen arm.
3. Update counts:

$$N_{A_t}(t) = N_{A_t}(t-1) + 1, \quad N_j(t) = N_j(t-1) \quad \text{for } j \neq A_t.$$

4. Update the empirical loss estimate for arm A_t via importance weighting:

$$\widehat{L}_{A_t}(t) = \frac{N_{A_t}(t-1) \widehat{L}_{A_t}(t-1) + K \ell_{A_t,t}}{N_{A_t}(t)},$$

and set $\widehat{L}_j(t) = \widehat{L}_j(t-1)$ for $j \neq A_t$.

Detailed Regret Analysis.

Setup: Assume K arms with losses $\ell_{i,t} \in [0, 1]$ and true mean μ_i for arm i . Under uniform sampling, when an arm is chosen its loss is scaled by K , so the empirical mean for arm i after n pulls is:

$$\widehat{\mu}_i(n) = \frac{1}{n} \sum_{s=1}^n K \ell_{i,t_s}.$$

Let $\mu^* = \min_i \mu_i$ and define the pseudo-regret as:

$$R_T = \sum_{t=1}^T \mathbb{E}[\ell_{A_t,t}] - T \mu^*.$$

Concentration Bound: Since each observed loss is in $[0, K]$, a concentration inequality (e.g., Hoeffding's) implies that with high probability,

$$|\widehat{\mu}_i(n) - \mu_i| \leq \sqrt{\frac{\alpha K^2 \ln T}{n}},$$

for some constant $\alpha > 0$.

UCB1 Index: Define the index for arm i at round t as:

$$I_i(t) = \widehat{\mu}_i(N_i(t)) + c \sqrt{\frac{\ln t}{N_i(t)}},$$

where c is chosen to cover the K -dependent variance. Standard analysis shows that the number of suboptimal pulls for any arm i (with gap $\Delta_i = \mu_i - \mu^*$) is bounded by

$$N_i(T) \leq \frac{C \ln T}{\Delta_i^2},$$

with C including factors from the importance weighting. Thus, the regret is bounded by

$$R_T \leq \sum_{i \neq i^*} \Delta_i N_i(T) = O\left(\sqrt{K T \ln T}\right)$$

in the adversarial or worst-case scenario.

Conclusion (a): The modified UCB1, with importance-weighted updates, yields a pseudo-regret of order

$$O\left(\sqrt{K T \ln T}\right),$$

up to constants depending on c and the loss distribution.

(b) Why the Modified UCB1 Cannot Exploit the Small Variance

In the modified UCB1, each observed loss is multiplied by K (since $p_i(t) = 1/K$). Even if the original loss distribution has small variance, the effective variance is inflated by the factor K . Since the algorithm builds confidence intervals based on worst-case deviations, the benefit of a small underlying variance is lost. In other words, the importance-weighted update makes the algorithm behave conservatively—as if the variance were higher—thus preventing exploitation of any small-variance structure.

(c) Modifying EXP3 for Importance-Weighted Losses (Uniform Sampling)

We consider a logging policy that selects arms uniformly at random (with probability $1/K$ each). The EXP3 algorithm is modified so that when arm A_t is selected and its loss $\ell_{A_t,t}$ is observed, an unbiased estimator for the loss of each arm is formed by

$$\tilde{\ell}_{i,t} = \begin{cases} K \ell_{A_t,t} & \text{if } i = A_t, \\ 0 & \text{otherwise.} \end{cases}$$

Then, the exponential weights update is applied:

$$w_i(t+1) = w_i(t) \exp(-\eta \tilde{\ell}_{i,t}),$$

with initial weights $w_i(1) = 1$ for all i .

Pseudo-Code for Modified EXP3:

1. **Initialization:** Set $w_i(1) = 1$ for each $i \in \{1, \dots, K\}$.
2. **For each round** $t = 1, 2, \dots, T$:
 - (a) **(Logging policy)** An arm A_t is chosen uniformly at random (i.e., with probability $1/K$).
 - (b) **(Observe loss)** Observe the loss $\ell_{A_t,t}$ for the chosen arm.
 - (c) **(Form importance-weighted loss)** For each arm i , define

$$\tilde{\ell}_{i,t} = \begin{cases} K \ell_{A_t,t} & \text{if } i = A_t, \\ 0 & \text{otherwise.} \end{cases}$$

- (d) **(Update weights)** Update

$$w_i(t+1) = w_i(t) \exp(-\eta \tilde{\ell}_{i,t}),$$

where $\eta > 0$ is the learning rate.

Detailed Regret Analysis.

Setup: Let $L_i = \sum_{t=1}^T \ell_{i,t}$ be the cumulative loss of arm i and $L^* = \min_i L_i$ be the loss of the best arm. The (pseudo-)regret is:

$$R_T = \sum_{t=1}^T \ell_{A_t,t} - L^*.$$

We aim to show that $\mathbb{E}[R_T] = O(\sqrt{KT \ln K})$.

Step 1: Master Inequality. Define the potential $\Phi_t = \ln W(t)$ where $W(t) = \sum_{i=1}^K w_i(t)$. By the exponential weights update and a standard inequality we have:

$$\Phi_{t+1} - \Phi_t \leq -\eta \sum_{i=1}^K p_t(i) \tilde{\ell}_{i,t} + \frac{\eta^2}{2} \sum_{i=1}^K p_t(i) \tilde{\ell}_{i,t}^2,$$

with $p_t(i) = \frac{w_i(t)}{W(t)}$. Summing over $t = 1$ to T ,

$$\Phi_{T+1} - \Phi_1 \leq -\eta \sum_{t=1}^T \sum_{i=1}^K p_t(i) \tilde{\ell}_{i,t} + \frac{\eta^2}{2} \sum_{t=1}^T \sum_{i=1}^K p_t(i) \tilde{\ell}_{i,t}^2.$$

Since $W(1) = K$ (so $\Phi_1 = \ln K$), we next lower-bound Φ_{T+1} .

Step 2: Lower Bound on Φ_{T+1} . For any fixed arm i ,

$$w_i(T+1) = \exp\left(-\eta \sum_{t=1}^T \tilde{\ell}_{i,t}\right).$$

Thus,

$$W(T+1) \geq w_{i^*}(T+1) = \exp\left(-\eta \sum_{t=1}^T \tilde{\ell}_{i^*,t}\right),$$

where i^* is the best arm. Taking logs,

$$\Phi_{T+1} \geq -\eta \sum_{t=1}^T \tilde{\ell}_{i^*,t}.$$

Step 3: Combine the Bounds. We have:

$$-\eta \sum_{t=1}^T \tilde{\ell}_{i^*,t} - \ln K \leq -\eta \sum_{t=1}^T \sum_{i=1}^K p_t(i) \tilde{\ell}_{i,t} + \frac{\eta^2}{2} \sum_{t=1}^T \sum_{i=1}^K p_t(i) \tilde{\ell}_{i,t}^2.$$

Rearrange to obtain:

$$\sum_{t=1}^T \sum_{i=1}^K p_t(i) \tilde{\ell}_{i,t} - \sum_{t=1}^T \tilde{\ell}_{i^*,t} \leq \frac{\ln K}{\eta} + \frac{\eta}{2} \sum_{t=1}^T \sum_{i=1}^K p_t(i) \tilde{\ell}_{i,t}^2.$$

Step 4: Relate to Actual Losses. Since only the chosen arm contributes,

$$\sum_{i=1}^K p_t(i) \tilde{\ell}_{i,t} = p_t(A_t) (K \ell_{A_t,t}).$$

Under uniform sampling, in expectation $p_t(A_t) = \frac{1}{K}$; hence,

$$\mathbb{E} \left[\sum_{i=1}^K p_t(i) \tilde{\ell}_{i,t} \right] = \ell_{A_t,t}.$$

Similarly, one shows that

$$\sum_{i=1}^K p_t(i) \tilde{\ell}_{i,t}^2 \leq K \ell_{A_t,t},$$

since $\ell_{A_t,t} \in [0, 1]$.

Taking expectations and summing over t ,

$$\mathbb{E} \left[\sum_{t=1}^T \ell_{A_t,t} \right] - \sum_{t=1}^T \ell_{i^*,t} \leq \frac{\ln K}{\eta} + \frac{\eta K}{2} \mathbb{E} \left[\sum_{t=1}^T \ell_{A_t,t} \right].$$

Since $\ell_{A_t,t} \leq 1$ implies $\mathbb{E}[\sum_{t=1}^T \ell_{A_t,t}] \leq T$, we have

$$\mathbb{E}[R_T] \leq \frac{\ln K}{\eta} + \frac{\eta K T}{2}.$$

Choosing

$$\eta = \sqrt{\frac{2 \ln K}{K T}},$$

balances the two terms and gives

$$\mathbb{E}[R_T] \leq 2\sqrt{2} \sqrt{K T \ln K}.$$

Thus, the expected regret is

$$\mathbb{E}[R_T] = O\left(\sqrt{K T \ln K}\right).$$

(d) Anytime Modification of EXP3

For the anytime version of EXP3 (which does not assume knowledge of a fixed horizon T), we replace the constant learning rate with one that depends on the current round t . A common choice is:

$$\eta_t = \sqrt{\frac{2 \ln K}{K t}}.$$

With this time-varying learning rate, the expected regret bound remains of the same order as the fixed-horizon version:

$$O\left(\sqrt{K T \ln K}\right),$$

though with a slightly larger constant factor than if T were known in advance.