

TME6 (RLD)

→ Finish TME 5 (Reinforce - actor-critic) !

→ PPO $\left[\begin{array}{l} \text{KL} \\ \text{clip} \end{array} \right]$

$$\mathcal{L}_{\Theta_r}(\theta) = \mathbb{E}_{\text{old}} \left(\frac{\pi_{\theta}(a_t | s_t)}{\pi_{\theta_{\text{old}}}(a_t | s_t)} \hat{A}_t \right)$$

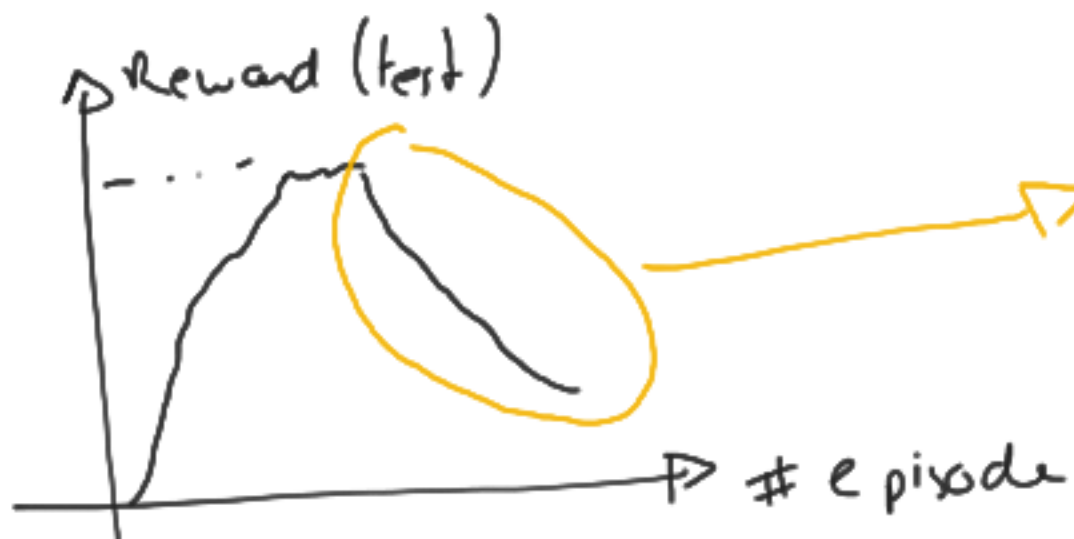
(KL) torch. distributions. Categorical (•)
dθ
probs = $\pi(a|s)$
logits = $\log \pi(a|s) + K$
torch. distributions. KL-divergence (dθ, dθ_{old})

(PPO) tensor. clamp(minimum, maximum) → $\begin{cases} x & \text{if } x \in [\text{min}, \text{max}] \\ \text{min} & \text{if } x < \text{min} \\ \text{max} & \text{if } x > \text{max} \end{cases}$

① Commencer par utiliser une baseline pour estimer l'avantage
 (→ marche par Cart Pole) → cf "réduction de la variance : baseline"

$$b(s_t) = \frac{1}{N} \sum_{\tau} R_t(\tau)$$

②



done = True lorsque l'épisode se termine
 ou est arrêté arbitrairement (≥ 500 pas de temps)

à ton tour

se termine

pas cared } R_t si done } $R_t + \gamma V_{\text{target}}^{\pi}(s_{t+1})$ sinon

correct } R_t si done \wedge not(truncated)
 $R_t + \gamma V_{\text{target}}^{\pi}(s_{t+1})$ sinon

↑
cible pour apprentissage V^{π}

ob, reward, done, info = env.step()

info.get("TimeLimit.truncated", False)

truncated =

True
si arrêt

False
sinon

$$\pi_{\Theta} \rightarrow \{(r_t, a_t, r_t, s_{t+1})\}$$

$$\text{ratio} = \underbrace{\begin{pmatrix} \pi_{\Theta}(s_1) \\ \vdots \\ \pi_{\Theta}(s_T) \end{pmatrix}}_{p(a_t | s_t, \Theta_{old})} \leftarrow \frac{\cdot}{\cdot}$$

$$\begin{pmatrix} \pi_{\Theta}(s_1) \\ \vdots \\ \pi_{\Theta}(s_T) \end{pmatrix}$$

$$P(a_t | s_t, \Theta)$$

↑
applies

Algorithm 4 PPO with Adaptive KL Penalty

Input: initial policy parameters θ_0 , initial KL penalty β_0 , target KL-divergence δ

for $k = 0, 1, 2, \dots$ **do**

Collect set of partial trajectories \mathcal{D}_k on policy $\pi_k = \pi(\theta_k)$

Estimate advantages $\hat{A}_t^{\pi_k}$ using any advantage estimation algorithm

Compute policy update

$$\theta_{k+1} = \arg \max_{\theta} \mathcal{L}_{\theta_k}(\theta) - \beta_k \bar{D}_{KL}(\theta || \theta_k)$$

by taking K steps of minibatch SGD (via Adam)

if $\bar{D}_{KL}(\theta_{k+1} || \theta_k) \geq 1.5\delta$ **then**

$$\beta_{k+1} = 2\beta_k$$

else if $\bar{D}_{KL}(\theta_{k+1} || \theta_k) \leq \delta/1.5$ **then**

$$\beta_{k+1} = \beta_k/2$$

end if

end for

calculate
 $\pi_{\theta_k}(a_t | s_t)$
 $p(a=0/s_t) \dots p(a=\dots/s_t)$
 \vdots
 $p(a=0/s_T) \dots p(a=\dots/s_T)$

utilize per step KL

actions = [...]

prob_a_s = probas [range(len(actions)), actions]

0	-	a ₀
1	-	a ₁
2	-	a ₂
⋮	⋮	
T-1	.	a _{T-1}

a while
gather