

# Lecture 8: PPO, DDPG / TD3

So far:



1) Perf. difference lemma:

$$\begin{aligned} J(\pi) - J(\mu) &= \mathbb{E}_{\tau \sim p_\pi(\tau)} \left[ \sum_{t=0}^{\infty} \gamma^t A^\pi(s_t, a_t) \right] \\ &= \frac{1}{1-\gamma} \mathbb{E}_{\substack{s \sim d^\pi(s) \\ a \sim \pi(\cdot | s)}} \left[ A^\pi(s, a) \right] \\ \max_{\pi} (J(\pi) - J(\mu)) \end{aligned}$$

2) Advantage weighted regression (AWR):

$$\max_{\pi} \mathbb{E}_{\substack{s \sim d^\pi(s) \\ a \sim \pi(\cdot | s)}} \left[ A^\mu(s, a) \right] - \alpha D(\pi, \mu)$$

Approximate by  $d^\mu(s)$  → this is ok if  $\pi$  &  $\mu$  are close in action dist.

Use  $D_{KL}(\pi(\cdot | s) \| \mu(\cdot | s))$  on states from  $\pi$

Gives us a nice closed form:

$$\pi^*(\cdot | s) \propto \mu(\cdot | s) e^{A(s, \cdot) / \alpha}$$

# PPO (Proximal policy optimization)

$$d^n \rightarrow d^m$$

$$\max_{\theta} \mathbb{E}_{s \sim d^m(s)} \left[ A^m(s, a) \right] - \alpha D_{KL}(\mu(\cdot|s) \| \pi_{\theta}(\cdot|s))$$

another way to

what if I want to learn compute tens is via IS on actions from  $\mu$  too??

$$\mathbb{E}_{x \sim p(x)} [f(x)] = \mathbb{E}_{x \sim q(x)} \left[ \frac{P(x)}{q(x)} f(x) \right]$$

different dist

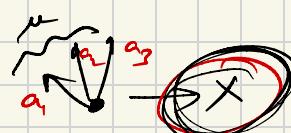
How can we do this for policy grad?

$$\nabla_{\theta} \mathbb{E}_{s \sim d^m(s)} \left[ \frac{\pi_{\theta}(a|s)}{\mu(a|s)} A^m(s, a) \right]$$

"Importance-weighted PG"

$$\mathbb{E} \left[ \frac{d^m(s)}{d(s)} \cdot A^m(s, a) \right]$$

$$\frac{\pi_\theta(a|s)}{\mu(a|s)} \cdot A^M(s, a)$$



Do I need the KL term?

Answer: No!

$$\max_{\theta} E_{\text{sample}} \left[ \frac{\pi_\theta(a|s)}{\mu(a|s)} \right] - D_{KL}(\mu || \pi_\theta)$$

- Can I just throw out the KL?

Not really as the update will be too large!

$$\Rightarrow \boxed{\pi_\theta \text{ is close to } \mu \text{ for all } s, a \Rightarrow \frac{\pi_\theta(a|s)}{\mu(a|s)} \text{ is close to 1}}$$

But  $D_{KL}(\mu || \pi)$  effectively keeps  $\pi$  &  $\mu$  close.

So, we can keep  $\frac{\pi_\theta(a|s)}{\mu(a|s)}$  close to 1.

$$\epsilon = 0.1, 0.2$$

hyper parameter

$$L_{clip}(\theta) = E_{\substack{s \sim d^m(s) \\ a \sim \mu(a|s)}} \left[ \underbrace{\text{clip} \left( \frac{\pi_\theta(a|s)}{\mu(a|s)}, 1-\epsilon, 1+\epsilon \right)}_{\text{---}} \cdot \underbrace{A^M}_{\text{---}} \right]$$

$$\frac{\pi_\theta(a|s)}{\mu(a|s)} > 1 + \epsilon$$

$$E_{s, a \sim \mu} \left[ \frac{\pi_\theta(a|s)}{\mu(a|s)} \cdot A^M(s, a) \right] \quad (1 - \epsilon) \cdot A^M$$

$$s, a \rightarrow \left| \frac{\pi_\theta(a|s)}{\mu(a|s)} \right| < 1 - \epsilon \Rightarrow \text{clip}(\dots) = \frac{1 - \epsilon}{1 + \epsilon}$$

$$r(\theta) = \tau_\theta(a|s) / \mu(a|s)$$

Any problems with the clip objective?

$$\Rightarrow \lambda^{\text{CLIP}(\theta)} = \mathbb{E} \left[ \exp \left( \underline{r(\theta)}, 1-\varepsilon, 1+\varepsilon \right) \cdot A^u \right]$$

- ① Can make an update even when we get clipped on the lower side!

$$d^{\text{new}}(\theta) = E \left[ \min \left( r(\theta) \cdot A^M, \text{clip}(r(\theta), 1-\epsilon, 1+\epsilon) \cdot A^U \right) \right]$$

when  $r(\theta) < 1 - \varepsilon$  &  $A^M(s, a) > 0$



action is now prob. under  $\pi_\theta$   
you see some signal.

but the  
action  
is a good  
action

but when  $r(0) < 1 - \varepsilon_{\text{bw}}$  &  $A^M(s, a) < 0$   
 you don't see signal.

instability

(Is this good or bad?)

Likewise if

highly likely under  $\pi$  but they are bad actions

$r(o) > 1 + \varepsilon$  but  $A^M(s, a) < 0$

you see some signal

"unclean" them

but when

$r(o) > 1 + \varepsilon$

but  $A^M(s, a) > 0$

there is no signal.

This hurts "exploration"!!

not in the RL case

↳ conservative choice

$\Rightarrow A^M$  can be erroneous

$\Rightarrow$  IS are high variance

sometimes useful to have Enough, Elow

"asymmetric clip"

used in DAPD (in LMs)

$$\begin{cases} 1 - \varepsilon \\ 1 + \varepsilon \end{cases}$$

$$[\Sigma_{\text{high}} \rightarrow \Sigma_{\text{low}}]$$

$$J(n) - J(m) = \mathbb{E}_{\pi} [A^n] - \mathbb{E}_{\pi} [A^m]$$

Option 2 for using the perf. difference lemma

use states from  $\mu$ , & compute  $A^\pi(s, a)$

$$Q^\pi(s, a)$$

can do this as

we can learn off policy



can be estimated using data from  
any policy

for any policy  $\pi$ ,

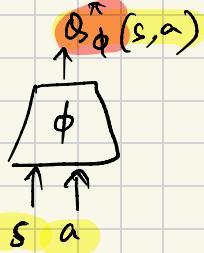
$$\forall s, a: Q^\pi(s, a) = r(s, a) + \gamma \mathbb{E}_{\substack{s' \sim p(\cdot | s, a) \\ a' \sim \pi(\cdot | s')}} [Q^\pi(s', a')]$$

"train a function"

We know how to implement this via ideas like

value-based methods.

$Q_\phi^\pi(s, a)$   
value-based RL



→ use target nets

→ soft/hard target updates

→ everything else remains the same

$$\pi_\theta(\cdot | s)$$

$$\max_{\theta} \mathbb{E}_\mu [Q_\phi^\pi(s, a)]$$

Loss for training the Q-function:

$$r(s, a) + \gamma \max_{a'} Q(s', a')$$

$$L(\phi) = \mathbb{E}_{\substack{s, a, r, s' \sim RBL \\ \text{Replay buffer}}} \left[ \left( Q_\phi(s, a) - \underbrace{y(s, a)}_{\text{!}} \right)^2 \right]$$

$$r(s, a) + \gamma \mathbb{E}_{\substack{s', a' \sim R \\ \text{tgt}}} \left[ Q_{\bar{\phi}}^{\text{tgt}}(s', a') \right]$$

How do we compute this?

$$\mathbb{E}_{\substack{s' \sim p(-|s, a) \\ a' \sim \pi(\cdot|s')}} \left[ Q_{\bar{\phi}}(s', a') \right] \approx$$

$\downarrow$

$$s' \in RB$$
$$a' \sim \pi_\theta(\cdot|s')$$

- 1) sample  $a'_0 \sim \pi_\theta(\cdot|s')$
- 2) query the tgt net. with it  $Q_{\bar{\phi}}(s', a')$

How do we do the policy update?

" $\pi_\theta$ "