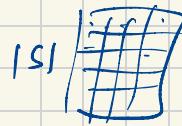


Lecture - 6 : Value-based RL

Last time :

① Value-iteration:

$$\pi(s) = \underset{a}{\operatorname{argmax}} Q(s, a)$$

$$Q(s, a) \leftarrow r(s, a) + \gamma V(s')$$


$$V(s) \leftarrow \max_{a'} Q(s, a')$$

② Fitted Q-iteration: (neural fQI)

{ for $i \geq 0, \dots$:

$$Q_{i+1} \leftarrow \underset{\theta}{\operatorname{argmin}} \mathbb{E}_{s, a, s'} \left[(Q_\theta(s, a) - y)^2 \right]$$

$$y(s, a) = r(s, a) + \gamma \max_{a'} \underline{Q_i(s, a')}$$

Q_i used here

$$y(s, a) = r(s, a) + \gamma \max_{a'} \underline{Q_i(s, a')}$$

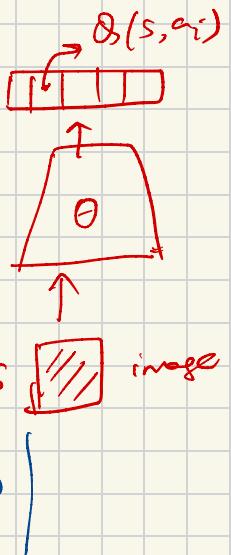
③ DQN:

- ↳ 1) replay buffer \rightarrow
- 2) target networks



history of
 (s, a, s')

Outline of DQN algorithm:



Exploration:

$$\pi_\theta(s)$$

$$\underset{a}{\operatorname{argmax}} \quad Q_{\theta_i}(s, a)$$

"exp-greedy" \rightarrow prob. ϵ

$$a \sim \text{unif}(a)$$

Update:



replay
buffer

$$(s, a, r, s')$$

prob $(1 - \epsilon)$, π_i

s image

only K steps of
grad. descent

$$E_{s, a, r, s' \sim \mathcal{D}}$$

Target networks

$$\textcircled{1} \quad \theta_0 = \theta_0^{\text{tgt}}$$

$$\theta_{i+1}$$

$$\max_a [Q_\theta - (\gamma \max_a Q_\theta^{\text{tgt}}(s', a'))]$$

$$\textcircled{2} \quad \theta_0 \rightarrow \text{grad. descent}$$

$$\text{main}$$

$$\text{target network}$$

update target net

$$\theta_K^{\text{tgt}} = \theta_K$$

$$\theta_i$$

(long)

UTD rate:

update-to-data

exploring
updates

Hard target update:

hard target update:



(few)
can't be
too large

vs soft

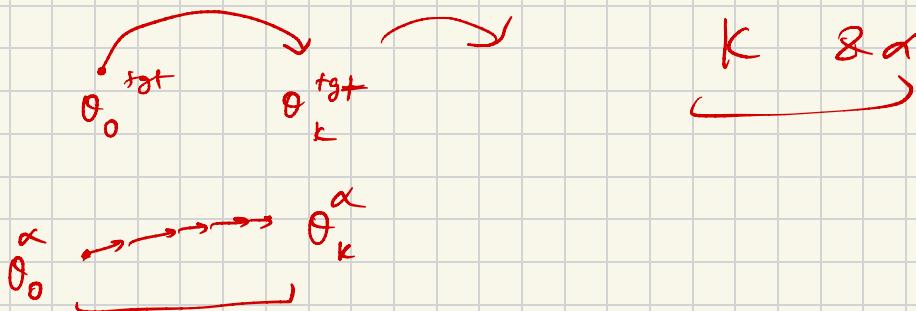
$$\alpha = 0.005$$

Soft target updates: $\left\{ \theta_{i+1}^{\text{tgt}} \leftarrow (1-\alpha) \theta_i^{\text{tgt}} + \alpha \theta_i \right.$
"soft"

$$\theta' \leftarrow (1-\alpha) \theta + \alpha \theta_{\text{new}}$$

hard $\left\{ \theta_k^{\text{tgt}} = \theta_k, \quad \theta_{2k}^{\text{tgt}} = \theta_{2k} \dots \dots \right.$

Why prefer one over the other?



Challenges with DQN:

"Double Q-learning"

① Overestimation in Q-values

error accumulation

$$y = r(s, a) + \gamma \max_{a'} Q_\theta(s', a')$$

tgt ✓

not trained on all s'

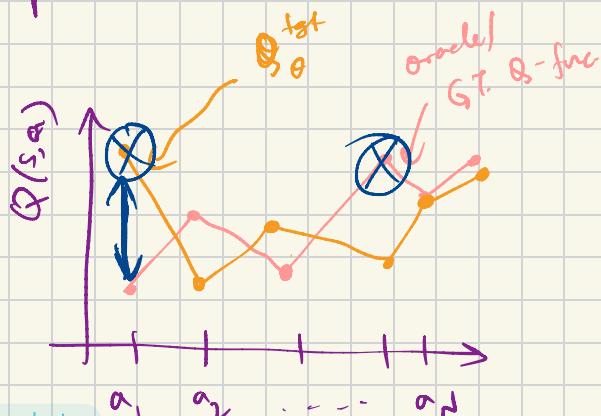
not trained on all a'

Compounding error

what is this term supposed to be?

optimal value from s' ideal $\boxed{V^*(s')}$ but
 Q is $Q_\theta \neq Q^*$

$\max_{a'} Q_\theta(s', a')$ → erroneous & we amplify the tve errors



Q_0^{tgt} → erroneous
 \Downarrow
 Q_0 → erroneous
 \Downarrow
 Q_1 → erroneous
 \Downarrow
 Q_K^{tgt} → erroneous

Double DQN

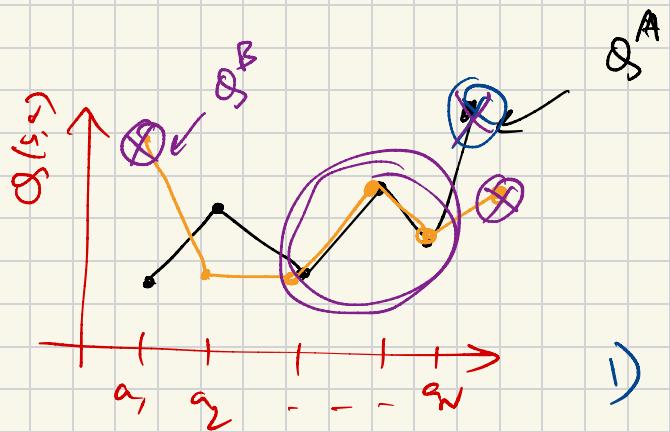
How can we fix this problem?

Idea: use argmax actions from one Q-function
 & compute Q-values under the other

$$\frac{(Q_{\theta}^A - y^A)^2}{Q_{\theta}^A} + \frac{(Q_{\theta}^B - y^B)^2}{Q_{\theta}^B}$$

$Q_{\theta}^{\text{tgt}, A}$ $Q_{\theta}^{\text{tgt}, B}$

initialized differently.



- 1) $\underset{a}{\operatorname{argmax}} Q_{\theta}(s, a)$
- 2) $Q_{\theta}^{\text{tgt}, A}(s, a^*)$

$$y = r(s, a) + \gamma \max_{a'} Q_{\theta}^{\text{tgt}}(s', a')$$

$$y^A = r(s, a) + \gamma Q_{\theta}^{\text{tgt}, A}(s', a'^*)$$

Another way to fix this problem.

Idea: reduce error compounding by using N-step returns.

$$\left(Q_0(s, a) - y(s, a) \right)^2$$

$$E_{s, a, s', r \sim \pi} [\dots]$$

$$\underline{y(s, a)} = \underline{r(s)} + \gamma \max_{\substack{\text{tgt} \\ a'}} Q(s', a')$$

\Downarrow

$$+ \arg \max_{a'} \quad \text{if seen in the replay buffer}$$

$$\underline{r(s, a)} + \gamma \underline{r(s', a')} + \gamma^2 \max_{a''} Q^{\text{tgt}}(s'', a'')$$

next state

$$\underline{Q(s, a)}$$

"N"

$$y(s, a) = r(s, a) + \gamma r(s_{n1}, a_{n1}) + \gamma^2 r(s_{n2}, a_{n2}) + \dots + \gamma^N \max_{a_{n+1}} Q^{\text{tgt}}(s_{n+1}, a_{n+1})$$

Practical details for Q-learning

- ① Coverage really helps.
- ② High UTD can destabilize training!
- ③ Improved loss functions:

$x \in \mathbb{R} \Rightarrow$  \Rightarrow Cross Entropy

Huber loss:

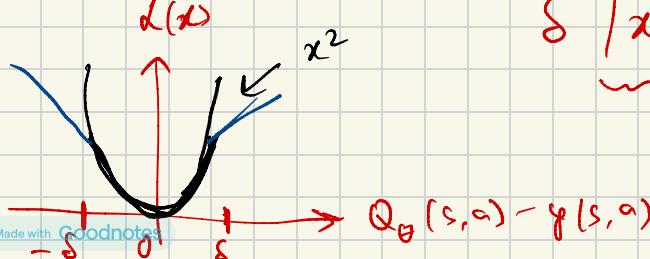
$$L(x) = \begin{cases} \frac{x^2}{2}, & |x| \leq \delta \\ \delta|x| - \frac{\delta^2}{2}, & |x| > \delta \end{cases}$$

$$x = (Q_\theta(s, a) - y(s, a))$$

$$(x)^2$$

$$\delta |x|$$

$$\frac{d |x|}{dx}$$



Advanced content:

Residual gradient vs TD-learning

$$\min_{\theta} \mathbb{E} \left[\| Q_{\theta}(s, a) - (r + \gamma Q_{\theta}(s', a')) \|_2^2 \right]$$

$$\Rightarrow \min_{\theta} \mathbb{E}_{s, a, s'} \left[(Q_{\theta}(s, a) - \gamma Q_{\theta}(s', a') - r(s, a))^2 \right]$$

What's the difference?

At convergence, in theory:

In practice:

Lecture 7: Advanced PG

So far, we saw:

Policy gradient methods
(policy, maybe an on policy $V(s)$)
on-policy

Q-learning methods
(no policy, just a value fn.)
off-policy

Can we build a hybrid of these methods?

Should be possible!

Why?

off-policy, actor-critic

off-policy

replay buffers

$\arg \max Q(s,a) \propto$

$\pi_\theta(a|s)$

How can we get started in a principled manner??

$$\mathbb{E}_{\pi} [A^{\mu}(s, a)] = 0$$

and

$$A^{\mu}(s, a)$$

any actions

Performance

Difference

Lemma

$$A^{\mu}(s, a) = Q^{\mu}(s, a) - V^{\mu}(s)$$

For any two policies $\pi(a|s)$ & $\mu(a|s)$, the following holds:

expected sum of discounted rewards under π

$$J(\pi) - J(\mu) = \mathbb{E}_{\tau \sim P_{\pi}(\tau)} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \right]$$

Return of policy π

policy μ

$$J(\pi) = \mathbb{E}_{\tau \sim P_{\pi}(\tau)} \left[\sum_t \gamma^t r(s_t, a_t) \right]$$

Proof:

$$\mathbb{E}_x [f(x)] = \mathbb{E}_{x_1, x_2, \dots} [f(x)]$$

$$J(\pi) = \mathbb{E}_{\tau \sim P_{\pi}(\tau)} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \right]$$

$$= \mathbb{E}_{s_0 \sim p(s_0)} [V^{\pi}(s_0)]$$

$$\begin{aligned} \tau &= \\ s_0 \sim p(s_0) &= \mathbb{E}_{\tau \sim P_{\pi}(\tau)} [V^{\pi}(s_0)] = \mathbb{E}_{\tau \sim P_{\mu}(\tau)} [V^{\pi}(s_0)] \quad \text{independent of whether it is in } \tau \text{ beyond } (s_0) \\ a_0 \sim \pi &= \mathbb{E}_{\tau \sim P_{\mu}(\tau)} \left[\sum_{t=0}^{\infty} \gamma^t V^{\pi}(s_t) - \sum_{t=1}^{\infty} \gamma^{t-1} V^{\pi}(s_{t-1}) \right] \\ s_1 \sim p(s_{t+1}|s_t, a_t) & \\ a_1 \sim \mu & \\ \vdots &= \mathbb{E}_{\tau \sim P_{\mu}(\tau)} \left[\sum_{t=0}^{\infty} \gamma^t (V^{\pi}(s_t) - \gamma V^{\pi}(s_{t+1})) \right] \end{aligned}$$

$$V^{\pi}(s_0) = V^{\pi}(s_0) + \gamma V^{\pi}(s_1) - \gamma V^{\pi}(s_1) + \gamma^2 V^{\pi}(s_2) - \gamma^2 V^{\pi}(s_2)$$

$$\underline{A^\pi(s_t, a_t)} = \underbrace{r(s_t, a_t) + \gamma}_{\text{red wavy line}} \cancel{v^\pi(s_{t+1}) - v^\pi(s_t)}$$

$$J(\pi) = \mathbb{E}_{\tau \sim p_{\mu}(\tau)} \left[\sum_{t=0}^{\infty} \gamma^t (v^\pi(s_t) - \gamma v^\pi(s_{t+1}) - r(s_t, a_t)) \right]$$

↓

$$J(\pi) = - \mathbb{E}_{\tau \sim p_{\mu}(\tau)} \left[\sum_{t=0}^{\infty} \gamma^t A^\pi(s_t, a_t) \right] + J(\mu)$$

↓

$$\Rightarrow J(\mu) - J(\pi) = \mathbb{E}_{\tau \sim p_{\mu}(\tau)} \left[\sum_{t=0}^{\infty} \gamma^t A^\pi(s_t, a_t) \right]$$

$$\max_{\pi} J(\pi) = \max_{\pi} (J(\pi) - J(\mu))$$

$$= \max_{\pi} \mathbb{E}_{s \sim d^\pi} [A^\mu(s, a)]$$

Why is this useful?

A practical off-policy policy gradient algorithm

$$J(\pi) - J(\mu) = \mathbb{E}_{\tau \sim P_\pi(\tau)} \left[\sum_{t=0}^{\infty} \gamma^t A^\mu(s_t, a_t) \right]$$

data collection policy

estimated via a critic

Option 1: use states from μ , A^μ instead

Option 2: use states from π , A^π instead

Which one should we use?

An approximation

$$\max_{\pi} J(\pi) - J(\mu)$$

$$E_{\tau \sim p_{\pi}(s)} \left[\sum_{t=0}^{\infty} \gamma^t A^{\mu}(s_t, a_t) \right]$$

$$= \sum_{t=0}^{\infty} E_{\substack{s_t \sim p_{\pi}(s_t) \\ a_t \sim \pi(\cdot | s_t)}} \left[\gamma^t A^{\mu}(s_t, a_t) \right]$$

ignore utilize

$$E_{\substack{s_t \sim p_{\mu}(s_t) \\ a_t \sim \pi(\cdot | s_t)}} \left[\gamma^t A^{\mu}(s_t, a_t) \right]$$

only term involved in maximization

Different ways to solve this maximization

when π & μ are close!!

Generic optimization problem:

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

$$-\alpha D(\pi, \mu)$$

what is this precisely?

We want $P_{\pi}(s_t) \& P_{\mu}(s_t)$ to be close!

$$P_{\pi}(s_t) - P_{\mu}(s_t) = ??$$

say π & μ are deterministic & agree with probability $1-\varepsilon$, then:

$$P_{\pi}(s_t) = (1-\varepsilon)^t P_{\mu}(s_t) + (1-(1-\varepsilon)^t) \text{Paiff}(s_t)$$

$$\Rightarrow |P_{\pi}(s_t) - P_{\mu}(s_t)| = (1 - (1-\varepsilon)^t) |P_{\mu}(s_t) - \text{Paiff}(s_t)| \leq 2 \cdot \varepsilon \cdot t$$

This means that:

$$D(\pi, \mu) = \underline{D_{TV}(\pi(-|s), \mu(-|s))}$$

total variations distance.

$$D_{TV}(P, Q) \leq \sqrt{\frac{1}{2} D_{KL}(P \| Q)}$$

}

$$\& D_{TV}(P, Q) \leq \sqrt{\frac{1}{2} D_{KL}(Q \| P)}$$

Pinsker's
inequality

Each of these choices gives a different algorithm:

① AWR $\rightarrow D_{KL}(\pi \| \mu)$

② PPO / TRPO $\rightarrow D_{KL}(\mu \| \pi)$

Advantage-weighted regression (AWR)

[Peng, Kumar, Liu, Levine, 2019]

$$\max_{\pi} \mathbb{E}_{s \sim d^{\mu}(s)} \left[\mathbb{E}_{a \sim \pi(\cdot | s)} [A^{\mu}(s, a)] - \alpha D_{KL}(\pi || \pi_0) \right]$$

Easier problem:

$$\max_{p(x)} \mathbb{E}_{x \sim p(x)} [f(x)] - \alpha D_{KL}(p || q)$$

$$\Rightarrow p^*(x) \propto q(x) \exp\left(\frac{f(x)}{\alpha}\right)$$

$$\text{So, } \pi^*(a|s) \propto \pi(a|s) \exp\left[\frac{A^{\mu}(s, a)}{\alpha}\right]$$

How do I use this oracle solution??

PPO (Proximal policy optimization)

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

another way to
compute this is via IS

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

How can we do this for policy grad?

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

Do I need the KL term?

Answer: No!

- Can I just throw out the KL?

Not really as the update will be too large!

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

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

$$\mathcal{L}^{\text{clip}}(\theta) = \mathbb{E}_{\substack{s \sim d^\mu(s) \\ a \sim \mu(a|s)}} \left[\text{clip}\left(\frac{\pi_\theta(a|s)}{\mu(a|s)}, 1-\epsilon, 1+\epsilon\right) A^u \right]$$

Any problems with the clip objective?

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

$$\mathcal{L}^{\text{CLIP}}(\theta) = \mathbb{E} \left[\min(r(\theta) \cdot A^u, \text{clip}(r(\theta), 1-\epsilon, 1+\epsilon) \cdot A^u) \right]$$

when $r(\theta) < 1-\epsilon$ & $A^u(s, a) > 0$
you see some signal.

but when $r(\theta) < 1-\epsilon$ & $A^u(s, a) < 0$
you don't see signal.

Is this good or bad?

Likewise if

$$r(0) > 1 + \varepsilon \quad \text{but} \quad A^M(s, a) < 0$$

you see some signal

$$\text{but when } r(0) > 1 + \varepsilon \quad \text{but} \quad A^M(s, a) > 0$$

there is no signal.

This hurts exploration !!

sometimes useful to have E_{high} , E_{low}

"asymmetric clip"

used in DAPD