

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

Option 2 for using the perf. difference lemma

use states from μ , & 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 π ,

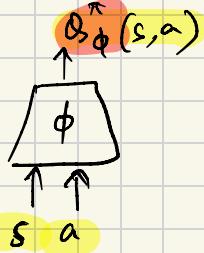
$$\forall s, a: Q^{\pi}(s, a) = r(s, a) + \gamma E_{\substack{s' \sim p(\cdot | s, a) \\ a' \sim \pi(\cdot | s')}} [Q^{\pi}(s', a')]$$

"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

$$\pi_{\phi}(\cdot | s)$$

$$\max_{\phi} E_{\mu} [Q_{\phi}^{\pi}(s, a)]$$

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?

" π_θ "

Policy update

We ideally want to do:

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

But there is a problem!!

How would you do gradient updates on π ?

$$\theta' \leftarrow \theta + \eta \nabla_{\theta} \mathbb{E}_{s, a \sim \pi} [A^{\pi}(s, a)]$$



$$A_{\phi}^{\pi}(s, a)$$

$$(\text{or } Q_{\phi}^{\pi}(s, a))$$

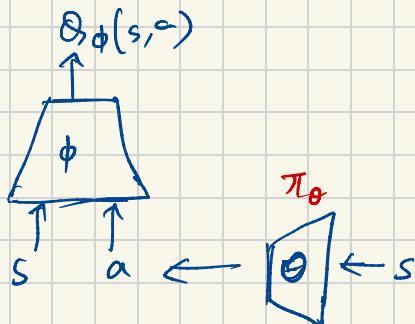
Another approximation

$$\theta' \leftarrow \theta + \eta \nabla_{\theta} \mathbb{E}_{s \sim d^{\pi}(s)} [Q_{\phi}^{\pi}(s, a)]$$
$$a \sim \pi_{\theta}(\cdot | s)$$

Back to something that looks like good, old PG!

But w/o the sum over time, which makes it easy!

How to compute? Just backprop through it!



$$\nabla_{\theta} Q_{\phi}^{\pi}(s, a) = \nabla_{\theta} Q_{\phi}^{\pi}(s, a_{\theta}(s))$$

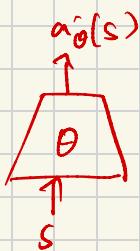
$$= \nabla_a Q_{\phi}^{\pi}(s, a) \Big|_{a=a_{\theta}(s)} \cdot \nabla_{\theta} a_{\theta}(s)$$

↓

What is this exactly?

1) say the policy is deterministic!

DDPG



2) Or it is stochastic!

SAC

$$\pi_{\theta}(\cdot | s) = \mathcal{N}(\mu_{\theta}(s), \Sigma_{\theta}(s))$$

Then $a \sim \pi_{\theta}(\cdot | s) \Leftrightarrow$ sample $\varepsilon \sim \mathcal{N}(0, I)$

$$\& a = \mu_{\theta}(s) + \varepsilon \cdot \sigma_{\theta}(s)$$

How to implement DDPG or its updated variant

TD3 correctly?

- 1) Two Q-networks, Two target Q-networks.
- 2) Deterministic policy: $\pi_\theta(s)$

Data collection: collect data using

$$a \sim \pi_\theta(s) + \varepsilon$$

↳ noise added to it

Training: sample $s, a, r, s' \sim$ Replay buffer.

$$y(s, a) = r(s, a) + \gamma \min \left(Q_{\phi, 1}^{\text{tgt}}(s', a'), Q_{\phi, 2}^{\text{tgt}}(s', a') \right)$$

$$a' \leftarrow \pi_\theta(s') + \varepsilon$$

Train $Q_{\phi, 1}$ & $Q_{\phi, 2}$. Train π_θ , update target networks.
once every d steps

Another algorithm - Soft Actor-Critic

✓
different from standard RL
parameterization

Standard RL: $\pi_{RL}^* = \operatorname{argmax}_{\pi} \mathbb{E}_{T \sim \pi} \left[\sum_t \gamma^t r(s_t, a_t) \right]$

Maxent RL:

$\pi_{\text{maxent}}^* = \operatorname{argmax}_{\pi} \mathbb{E}_{T \sim \pi} \left[\sum_{t=0}^{\infty} \gamma^t (r(s_t, a_t) + H(\pi)) \right]$

We can simplify the maxent RL objective!

$$H(\pi(\cdot | s_t)) = \mathbb{E}_{a \sim \pi(\cdot | s_t)} [-\log \pi(a | s_t)]$$

$$\pi_{\text{maxent}}^* = \max_{\pi} \mathbb{E}_{\tau \sim \pi} \left[\sum_{t=0}^{\infty} \gamma^t (r(s_t, a_t) + \alpha \mathbb{E}_{a_t \sim \pi} [-\log \pi]) \right]$$

Can move this here!!

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

"backup the entropy"

modified reward

$$\bar{r}_\alpha(s, a) = r(s, a) - \alpha \log \pi(a | s)$$

SAC := TD3 + Max-ent RL
(roughly)
very roughly !!

Q-functrion:

$$Q(s, a) \leftarrow r(s, a) + \gamma \mathbb{E}_{a' \sim \pi} [Q(s', a')]$$



$$Q(s, a) \leftarrow r(s, a) + \gamma \mathbb{E}_{a' \sim \pi} [Q(s', a') - \alpha \log \pi(a'|s)]$$

Policy extraction:

$$\pi_\theta$$

$$\theta' \leftarrow \theta + \gamma \nabla_\theta \mathbb{E}_{a \sim \pi_\theta^\phi(\cdot | s)} [Q_\phi(s, a)]$$



$$\theta' \leftarrow \theta + \gamma \nabla_\theta \mathbb{E}_{a \sim \pi_\theta^\phi(\cdot | s)} [Q_\phi(s, a) - \alpha \log \pi_\theta(a | s)]$$

Why use entropy in here?

Practical design choices

- ① Use two critic networks

$$\min (\mathcal{Q}_1, \mathcal{Q}_2)$$

- ② No noise during exploration. Just sample from the stochastic policy

- ③ α is dynamically tuned!

Why does adding the entropy help?

→ Stochastic policy is easier to optimize

→ More ways to succeed !!