

Back to J

- Objective:

$$\begin{aligned} J(\theta) &= \sum_s \mu_{\pi_\theta}(s) v_{\pi_\theta}(s) \\ &= \sum_s \mu_{\pi_\theta}(s) \sum_a \pi_\theta(a|s) q_{\pi_\theta}(s, a) \end{aligned}$$

- True gradient:

$$\begin{aligned} \nabla J(\theta) &\propto \sum_s \mu_{\pi_\theta}(s) \sum_a \nabla \pi_\theta(a|s) q_{\pi_\theta}(s, a) \\ &\propto \sum_s \mu_{\pi_\theta}(s) \sum_a \nabla \pi_\theta(a|s) (q_{\pi_\theta}(s, a) - v_{\pi_\theta}(s)) \end{aligned}$$

- Stochastic gradient:

$$\begin{aligned} \widehat{\nabla} J(\theta) &\propto \nabla \log \pi_\theta(A_t|S_t) (q_{\pi_\theta}(S_T, A_T) - v_{\pi_\theta}(S_t)) \\ &\propto \nabla \log \pi_\theta(A_t|S_t) a_{\pi_\theta}(S_T, A_T) \end{aligned}$$

- On policy algorithm if we can estimate $a_{\pi_\theta}(S_T, A_T) = q_{\pi_\theta}(S_T, A_T) - v_{\pi_\theta}(S_t)$.
(Critic)

Off-Policy J

- Objective:

$$\begin{aligned} J_b(\theta) &= \sum_s \mu_b(s) v_{\pi_\theta}(s) \\ &= \sum_s \mu_b(s) \sum_a \pi_\theta(a|s) q_{\pi_\theta}(s, a) \end{aligned}$$

- True gradient:

$$\nabla J_b(\theta) = \sum_s \mu_b(s) \sum_a (\nabla \pi_\theta(a|s) q_{\pi_\theta}(s, a) + \pi_\theta(a|s) \nabla q_{\pi_\theta}(s, a))$$

- ∇q_{π_θ} term hard to compute!
- Descent direction:

$$\tilde{\nabla} J_b(\theta) = \sum_s \mu_b(s) \sum_a \nabla \pi_\theta(a|s) q_{\pi_\theta}(s, a) = \sum_s \mu_b(s) \sum_a \nabla \pi_\theta(a|s) a_{\pi_\theta}(s, a)$$

- Stochastic descent direction

$$\hat{\tilde{\nabla}} J_b(\theta) = \frac{\pi_\theta(A_t|S_t)}{b(A_t|S_t)} \nabla \log \pi_\theta(a|s) a_{\pi_\theta}(S_t, A_t)$$

- Off-policy algorithm if we can estimate $a_{\pi_\theta}(S_t, A_t)$ (Critic)

Trust Region Policy Optimization

- Local objective:

$$J_{\pi_{\theta_{old}}}(\theta) = \sum_s \mu_{\pi_{\theta_{old}}}(s) \sum_a \pi_{\theta}(a|s) a_{\pi_{\theta_{old}}}(s, a)$$

- If convergence we recover the on-policy goal.

- True gradient:

$$\nabla J_{\pi_{\theta_{old}}}(\theta) = \sum_s \mu_{\pi_{\theta_{old}}}(s) \sum_a \nabla \pi_{\theta}(a|s) a_{\pi_{\theta_{old}}}(s, a)$$

- Identical to the descent direction of the off-policy algorithm at θ_{old} .
- Optimization of the local objective only in a neighborhood of π_{old} :

$$\sup_s \text{KL}(\pi_{\theta_{old}}(s), \pi_{\theta}(s)) \leq \epsilon$$

- Strong link with a backtracking algorithm of the off-line version.
- Need to replace the trust region by an approximate one based on

$$\mathbb{E}_{\pi_{\theta_{old}}} [\text{KL}(\pi_{\theta_{old}}(s), \pi_{\theta}(s))] \quad (\text{quadratic constraint in } \theta)$$

Proximal Policy Optimization

- State Of The Art actor-critic algorithm.

- First version:

$$J_{\pi_{\theta_{old}}}(\theta) = \sum_s \mu_{\pi_{old}}(s) \sum_a \pi_{\theta}(a|s) a_{\pi_{\theta_{old}}}(s, a) + \lambda \mathbb{E}_{\pi_{\theta_{old}}} [\text{KL}(\pi_{\theta_{old}}(s), \pi_{\theta}(s))]$$

- Very similar to TRPO.

$$= \sum_s \mu_{\pi_{old}} \sum_a \pi_{old} \frac{\pi_{\theta}}{\pi_{old}} \dots - \epsilon \text{KLO}$$

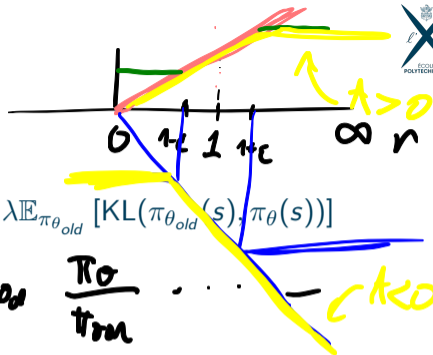
- Second version:

$$J_{\pi_{\theta_{old}}}(\theta) = \sum_s \mu_{\pi_{old}}(s) \sum_a \pi_{\theta_{old}}(a|s)$$

$$\times \min \left(r_{\theta_{old}, \theta}(a, s) a_{\pi_{\theta_{old}}}(s, a), \text{clip}(r_{\theta_{old}, \theta}(a, s), 1 - \epsilon, 1 + \epsilon) a_{\pi_{\theta_{old}}}(s, a) \right)$$

with $r_{\theta_{old}}(\theta) = \frac{\pi_{\theta}(a|s)}{\pi_{\theta_{old}}(a|s)}$.

- Worst case scenario for the advantage as soon as the ratio is away from 1.
- No theoretical justification but simple and efficient algorithm.



- Modification of the reward to favor high entropy policy:

$$R_t \rightarrow R_t + \lambda \mathcal{H}(\pi(S_t))$$

- Classical state value function:

$$v(s) = \mathbb{E} \left[\sum_{i=0}^{\infty} \gamma^i (R_{t+i+1} + \lambda \mathcal{H}(\pi(\cdot | S_{t+i}))) \mid S_t = s \right]$$

- Modified state action value function defined by

$$v(s) = \sum_a \pi(a|s) (q(s, a) - \lambda \log(\pi(a|s)))$$

- Fixed point operator:

$$\begin{aligned} \mathcal{T}^\pi q(s, a) &= r(s, a) + \mathbb{E} [\gamma v(s', a)] \\ &= r(s, a) + \mathbb{E}_\pi [q(s', a') - \lambda \log(\pi(a'|s'))] \end{aligned}$$

$\pi^* = \operatorname{argmax}_\pi q(\cdot, \cdot)$
 \uparrow
 $\pi \propto e^{-\frac{1}{\lambda} q(s, a)}$

- Policy improvement rule:

$$\pi^+(\cdot|s) = \operatorname{argmax}_{\pi(\cdot|s)} \sum_a \pi(a|s) (q(s, a) - \lambda \log(\pi(a|s)))$$

$$\pi^+(a|s) \propto \exp\left(-\frac{1}{\lambda} q(s, a)\right)$$

implies $v_{\pi^+}(s) \geq v_{\pi}(s)$.

- Stronger link between the critic and the actor with an ideal update following the policy improvement rule.

- Parametric π_θ , Q_ϕ and $Q_{\phi'}$...

- Optimization in ϕ :

$$V_{\phi'}(s, a) = \mathbb{E}_{\pi_\theta} [Q_{\phi'}(S_t, A_t) - \lambda \log \pi_\theta(A_t|S_t)]$$

$$J(\phi) = \mathbb{E}_{\pi_\theta} \left[(Q_\phi(S_t, A_t) - (r(S_t, A_t) + \gamma \mathbb{E} [V_{\phi'}(S_{t+1})]))^2 \right]$$

$$\hat{\nabla} J(\phi) = 2 \nabla Q(S_t, A_t)$$

$$\times (Q_\phi(S_t, A_t) - (r(S_t, A_t) + \gamma (Q_{\phi'}(S_{t+1}, A_{t+1}) - \lambda \log \pi_\theta(A_{t+1}|S_{t+1}))))$$

- ϕ' : slow version of ϕ (exponential average)

$$\phi' \leftarrow (1-\rho) \phi' + \rho \phi$$

- Two-scales trick!

- Optimization in θ :

$$\begin{aligned} J(\theta) &= \mathbb{E}_{\pi_{\theta}} \left[\text{KL}(\pi_{\theta}(\cdot|S_t), e^{\frac{1}{\lambda} Q_{\phi}(S_t, \cdot)} / Z(S_t, \phi)) \right] \\ &= \mathbb{E}_{\pi_{\theta}} \left[\sum_a \pi_{\theta}(a|S_t) \left(-\log \pi_{\theta}(a|S_t) + \frac{1}{\lambda} Q_{\phi}(S_t, a) \right) \right] + \text{Cst} \end{aligned}$$

- Easy optimization when neglecting the effect of π_{θ} in $S_t \dots$

$$\hat{\nabla} J(\theta) = \sum_a \left(-\log \pi_{\theta}(a|S_t) + \frac{1}{\lambda} Q_{\phi}(S_t, a) - 1 \right) \nabla_{\theta} \pi_{\theta}(a|S_t)$$

- θ is updated at a much slower pace than ϕ .
- Two-scales algorithm trick again!

$$a = \mu + \sigma \epsilon$$

- Adaptation possible to continuous action using the reparametrization trick:

$$A_t = \Phi_\theta(S_t, \epsilon_t)$$

with ϵ_t a known density $p(\epsilon)$ and Φ an invertible transform ($\epsilon_t = \Phi_\theta^{-1}(S_t, A_t)$)

- Implicit parametrization of π_θ :

$$\pi_\theta(a|s) = |J_{\Phi_\theta^{-1}(s,a)}| p(\Phi_\theta^{-1}(s, a))$$

- Rewriting of the objective:

$$\begin{aligned} J(\theta) &= \mathbb{E}_{\pi_\theta} \left[\text{KL}(\pi_\theta(\cdot|S_t), e^{\frac{1}{\lambda} Q_\phi(S_t, \cdot)} / Z(S_t, \theta)) \right] \\ &= \mathbb{E}_{\pi_\theta, \epsilon} \left[-\log \pi_\theta(\Phi_\theta(S_t, \epsilon)|S_t) + \frac{1}{\lambda} Q_\phi(S_t, \Phi_\theta(S_t, \epsilon)) \right] + Cst \end{aligned}$$

- Easy optimization when neglecting the effect of π_θ in $S_t \dots$

$$\begin{aligned} \widehat{\nabla} J(\theta) &= -\nabla_\theta \log \pi_\theta(\Phi_\theta(S_t, \epsilon)|S_t) \\ &\quad + \nabla_\theta \Phi_\theta(S_t, \Phi_\theta(S_t, \epsilon)) \left(-\nabla_a \log \pi_\theta(\Phi_\theta(S_t, \epsilon)|S_t) + \frac{1}{\lambda} \nabla_a Q_\phi(S_t, \Phi_\theta(S_t, \epsilon)) \right) \end{aligned}$$

Double Q learning and extension

- Classical Q learning:
 - Target: $R_{s,a} + \max_{a'} Q_{\phi}(s', a')$
 - Approximation issue: $Q_{\phi}(s', a') \sim Q(s, a) + \epsilon(s, a)$
 - Consequence: $\mathbb{E} [\max_a Q_{\phi}(S_t, a)] \geq \max(Q(s, a) + \mathbb{E} [\epsilon(s, a)])$
- Double Q learning:
 - Two Q learning function: $Q_{\phi_i}(s, a)$
 - Used in a crossed way for the target of Q_{ϕ_i} :
$$R_{s,a} + Q_{\phi_{i'}}(s', \underset{a'}{\operatorname{argmax}} Q_{\phi_i}(s', a'))$$
 - Mitigate the bias.
- Similar overestimation bias issue in actor critic approach.
- Clipped double Q learning:
 - Two Q learning function: $Q_{\phi_i}(s, a)$
 - Used in a pessimistic way for the target of Q_{ϕ_i} :
$$R_{s,a} + \min_i Q_{\phi_i'}(s', a') - \lambda \log \pi_{\theta}(a'|s')$$as well as in the optimization on θ .

- In most of the algorithm, the expectation along trajectory is replaced by an empirical average over past short pieces of trajectories stored in a replay buffer.
- Corresponds exactly to the idea of an empirical simulator when the policy is fixed.
- If the policy is changing across time, we should use a importance sampling correction to be faithful with the theory. . .
- Not necessary for one-step Q learning but necessary for the actor-critic approach as the stationary law on the state is used.
- Not an issue in practice: *neglecting the effect of π_θ in S_t* of the previous slides.
- To even reduce the issue: use only *recent* trajectories in the buffer.