

SOFT ACTOR-CRITIC

- OFF-POLICY, reuse past experience (not directly feasible with standard policy gradient)
- MAXIMUM ENTROPY
- STOCHASTIC ACTION

→ ON-POLICY requires new samples & gradient step
 → augments the standard maximum reward RL objective with an ENTROPY MAXIMIZATION TERM
 ↳ better exploration and robustness in estimation error
 ↳ same prob. of acting for example two actions that have the same (maximum) entropy

3 ingredients:

- 1) Actor-Critic architecture with separate policy and value fn
- 2) OFF-POLICY to reuse previously collected data
- 3) Entropy maximization to enable stability and exploration

1) Start from POLICY ITERATION $\left\{ \begin{array}{l} \rightarrow \text{POLICY EVALUATION} \\ \rightarrow \text{POLICY IMPROVEMENT using the value fn} \end{array} \right.$
 Policy \rightarrow Actor
 Value fn \rightarrow Critic

3) Maximum entropy RL optimizes policies to maximize both the expected return and the expected entropy of the policy
 ↳ continuous state and action spaces

$$\text{MDP} = (\underbrace{S, A, p, r}_{\text{continuous state and action spaces}}, \underbrace{r_{min}, r_{max}}_{\text{range of rewards}}) \rightarrow S \times A \rightarrow [r_{min}, r_{max}]$$

$\rightarrow S \times S \times A \rightarrow [0, \infty) = \text{prob of } S_{t+1} \in S \text{ given } S_t \in S \text{ and } A_t \in A$

$$\begin{aligned} \pi(\cdot | s_t) &= \text{state} \\ \pi(\cdot | s_t, a_t) &= \text{state-action} \end{aligned} \quad \left. \begin{array}{l} \text{marginals of trajectory dist. induced by } \pi(a_t | s_t) \\ \text{of EM term} \end{array} \right\}$$

$$J(\pi) = \sum_{t=0}^T \mathbb{E}_{(s_t, a_t) \sim P^\pi} \underbrace{\left[r(s_t, a_t) + \alpha H(\pi(\cdot | s_t)) \right]}_{\text{STANDARD reward}} \quad (1)$$

$\xrightarrow{\text{TEMPERATURE parameter determines the importance}}$

- wider exploration \rightarrow gives up on clearly unpromising avenues
- can capture multiple modes of near-optimal behavior
 \hookrightarrow equal prob. to those actions

Entropy Measure of chaos in terms of "how random a RV is"

$$\rightarrow \underline{H(P)} = \mathbb{E}_{x \sim P} [-\log P(x)], \times \text{RV with distribution } P$$

SOFT POLICY ITERATION 1) SOFT POLICY, compute value of π according to ME objective in (1)

For a fixed π , the SOFT Q-value can be computed iteratively from any $Q: S \times A \rightarrow \mathbb{R}$ and a modified Bellman operator T^π

$$\star T^\pi Q(s_t, a_t) \triangleq r(s_t, a_t) + \gamma \mathbb{E}_{s_{t+1} \sim P} [V(s_{t+1})],$$

$$\text{where } V(s_t) = \mathbb{E}_{a_t \sim \pi} \left[Q(s_t, a_t) - \underbrace{\log \pi(a_t | s_t)}_{\text{ENTROPY!}} \right]$$

SOFT STATE-VALUE FN

LEMMA Q^k (with $Q^{k+1} = \pi^\pi Q^k$) will converge to the soft Q-values of π as $k \rightarrow \infty$

2) SOFT POLICY, update policy towards the exp of the new Q-fn IMPROVEMENT

$$\star \quad \pi_{\text{new}} = \arg \min_{\pi' \in \Pi} D_{KL} \left(\pi'(\cdot | s_t) \parallel \frac{\exp(Q_{\text{old}}^\pi(s_t, \cdot))}{Z_{\text{old}}^\pi(s_t)} \right) \quad (2)$$

we restrict Π to
a Gaussian family

↓ convenient to
use KULLBACK - LEIBLER
DIVERGENCE

↓ PARTITION FN that normalizes
the distib.

Can be ignored because it is
intractable and does not
effect π_{new}

LEMMA $Q^{\text{new}}(s_t, a_t) \geq Q^{\text{old}}(s_t, a_t)$, $\forall (s_t, a_t) \in S \times A$ and $\pi_{\text{old}}, \pi_{\text{new}} \in \Pi$

But we need to approximate it to represent the Q-values
(too expensive and work only in the tabular case)

⇒ SOFT ACTOR-CRITIC to deal with action/state spaces that are continuous

We use a parametrized $V_\phi(s_t)$, $Q_\theta(s_t, a_t)$ and $\pi_\phi(a_t | s_t)$ where ψ , θ and ϕ are the params of the network

↓ value fn models as NNs, policy as Gaussian with μ and σ given by NNs

→ Soft value fn trained to minimize the SQUARED RESIDUAL ERROR

$$\star J_V(\psi) = \mathbb{E}_{D^{\text{END}}} \left[\frac{1}{2} \left(V_\psi(s_t) - \mathbb{E}_{a_t \sim \pi_\phi} [Q_\theta(s_t, a_t) - \log \pi_\phi(a_t | s_t)] \right)^2 \right]$$

where D is the distib. of previously sampled states and actions
↓ REPLAY BUFFER

→ Soft Q-fn params can be trained to minimize the soft Bellman residual

$$\star J_Q(\theta) = \mathbb{E}_{(s_t, a_t) \sim D} \left[\frac{1}{2} \left(Q_\theta(s_t, a_t) - \hat{Q}(s_t, a_t) \right)^2 \right] \quad \text{as usual} \leftarrow (r(s_t, a_t) + \gamma \mathbb{E}_{s_{t+1} \sim P} [V_\phi(s_{t+1})])$$

N.B.: Both $J_V(\psi)$ and $J_Q(\theta)$ can be optimized with stochastic gradients

→ Finally, policy params learned by minimizing the expected D_{KL} in (2)

But no way to back-propagate with the sampling

⇒ RE PARAMETERIZATION TRICK using a NN transformation, so samples are determined according to

$$\star \tilde{\delta}_\theta(s, \{\}) = \tanh(\mu_\theta(s) + \sigma_\theta(s) \odot \{\}), \{\} \sim N(\emptyset, I) \quad (\text{SPHERICAL GAUSSIAN})$$

Now, gradients flow through μ_θ and σ_θ (it's the only- way to allow it)

Two Q-fns to mitigate positive bias in policy improvement step that degrades performance of value-based method

↓ these Q-fns are parametrized with θ_i and train independently to optimize $J_Q(\theta_i)$

↓ then use the minimum of the Q-fns for the value gradient