

(soft, smooth, stable)
SAC: 引入最大熵, 使 action 探索能力加强, 在 PPO 里, 当在线和离线策略的 KL 最小, 则取样一致, SAC 是学习一个分布, 每条数据的不确定性 $= -\log \pi(a|s)$, 尽量造 Entropy 大的 batch

```
def actor_learn(self, obs):
    action, log_pi = self.sample_obs(
        qf1_pi, qf2_pi = self.critic.value(obs)
    )
    min_qf_pi = layers.elementwise_min(qf1_pi, qf2_pi)
    cost = log_pi * self.alpha - min_qf_pi # 0.2, 0.2
    cost = layers.reduce_mean(cost)
    optimizer = fluid.optimizer.AdamOptimizer(self.actor_lr)
    optimizer.minimize(cost, parameter_list=self.actor_parameters())

    return cost
```

在 TD3 里, actor cost = $Q(s, a)$, 这增加了 entropy

$$J_{\pi}(\phi) = \mathbb{E}_{a_t \sim \pi} \left[D_{KL} \left(\pi_{\phi}(\cdot | s_t) \parallel \frac{\exp(Q_{\phi}(s_t, \cdot))}{Z_{\phi}(s_t)} \right) \right]$$

$$J_{\pi}(\phi) = \mathbb{E}_{a_t \sim \pi} \left[-\log \pi_{\phi}(a_t | s_t) - Q_{\phi}(s_t, a_t) \right]$$

$$\nabla_{\phi} J_{\pi}(\phi) = \nabla_{\phi} \log \pi_{\phi}(a_t | s_t) + (\nabla_{a_t} \log \pi_{\phi}(a_t | s_t) - \nabla_{a_t} Q_{\phi}(s_t, a_t)) \nabla_{\phi} Q_{\phi}(s_t, a_t)$$

Algorithm 1 Soft Actor-Critic $A(s) = \text{mean}, \log \pi(a|s) \in (-20, 2)$

Input: θ_1, θ_2, ϕ
 $\theta_1 \leftarrow \theta_1, \theta_2 \leftarrow \theta_2$
 $D \leftarrow \emptyset$

for each environment step do
 for each gradient step do
 $a_t \sim \pi_{\phi}(a_t | s_t)$
 $s_{t+1} \sim p(s_{t+1} | s_t, a_t)$
 $D \leftarrow D \cup \{(s_t, a_t, r(s_t, a_t), s_{t+1})\}$
end for
for each gradient step do
 $\theta_1 \leftarrow \theta_1 - \lambda_Q \nabla_{\theta_1} J_Q(\theta_1)$ for $i \in \{1, 2\}$
 $\theta_2 \leftarrow \theta_2 - \lambda_V \nabla_{\theta_2} J_V(\theta_2)$
 $\alpha \leftarrow \alpha - \lambda \nabla_{\alpha} J(\alpha)$
 $\theta_1 \leftarrow \tau \theta_1 + (1 - \tau) \theta_1$ for $i \in \{1, 2\}$
end for
Output: θ_1, θ_2, ϕ

Initial parameters
 Update policy weights
 Initialize target network weights
 Initialize an empty replay pool
 Sample action from the policy
 Sample transition from the environment
 Store the transition in the replay pool
 Update the Q-function parameters
 Update policy weights
 Adjust temperature
 Update target network weights
 Optimized parameters

$$y = r + \gamma \cdot \min_a Q_{\phi}(s', a) - \log \pi$$

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

$$\nabla_{\psi} J_V(\psi) = \nabla_{\psi} V_{\psi}(s_t) (V_{\psi}(s_t) - Q_{\phi}(s_t, a_t) + \log \pi_{\phi}(a_t | s_t))$$

行为策略: $\text{alg. sample}(s)$ 左图
 评价策略: $\tanh(A(s)) \cdot a_{\max} \leftarrow ??$
 (模型) (更新)

```
def critic_learn(self, obs, action, reward, next_obs, terminal):
    next_obs_action, next_obs_log_pi = self.sample(next_obs)
    qf1_next_target, qf2_next_target = self.target_critic.value(
        next_obs, next_obs_action
    )
    min_qf_next_target = layers.elementwise_min(
        qf1_next_target, qf2_next_target
    )
    next_obs_log_pi = self.alpha * next_obs_log_pi + self.alpha

    terminal = layers.cast(terminal, dtype='float32')
    target_Q = reward + (1.0 - terminal) * self.gamma * min_qf_next_target
    target_Q = layers.square_error_cost(current_Q1, target_Q) + layers.square_error_cost(current_Q2, target_Q)
    cost = layers.reduce_mean(cost)
    optimizer = fluid.optimizer.AdamOptimizer(self.critic_lr)
    optimizer.minimize(cost)
    return cost
```

DDPG

```
def critic_learn(self, obs, action, reward, next_obs, terminal):
    next_obs_action = self.target_model.policy(next_obs)
    next_Q = self.target_model.value(next_obs, next_obs_action)

    terminal = layers.cast(terminal, dtype='float32')
    target_Q = reward + (1.0 - terminal) * self.gamma * next_Q
    target_Q.stop_gradient = True

    Q = self.model.value(obs, action)
    cost = layers.square_error_cost(Q, target_Q)
    cost = layers.reduce_mean(cost)
    optimizer = fluid.optimizer.AdamOptimizer(self.critic_lr)
    optimizer.minimize(cost)
    return cost
```

前 3 万步, actor-lr: $1e-3$ critic-lr: $1e-3$ γ : 0.99 tau: 0.9
 noise: 0.25 batch size: 256
 后面: actor-lr: $5e-3$ critic-lr: $5e-3$ γ : 0.999 tau: 0.01
 noise: 1 batch size: 256
 不再解与面为啥动作
 噪音反而增加?
 fc1: 100 $A(s) \in (-1, 1)$ 探索+集
 加噪: $a = \text{np.clip}(\text{gaussian}(a, \text{noise}), -1, 1)$ 这噪音多大
 For DDPG: $a' = Q_{\phi}(s')$ $y = r + \gamma \cdot V_{\phi}(s'; a')$

TD3 算法 (3 次更新)

```
def critic_learn(self, obs, action, reward, next_obs, terminal):
    noise = layers.gaussian_random_batch_size_like(
        action, shape=[1, action.shape[1]]
    )
    noise = layers.clip(
        noise * self.policy_noise, # 0.2
        min=self.noise_clip, # 0.5
        max=self.noise_clip, # 0.5
    )
    next_obs_action = self.target_model.policy(next_obs) + noise
    next_obs_action = layers.clip(next_obs_action, -self.max_action, self.max_action)

    next_Q1, next_Q2 = self.target_model.value(next_obs, next_obs_action)
    next_Q = layers.elementwise_min(next_Q1, next_Q2)

    terminal = layers.cast(terminal, dtype='float32')
    target_Q = reward + (1.0 - terminal) * self.gamma * next_Q
    target_Q.stop_gradient = True

    current_Q1, current_Q2 = self.model.value(obs, action)
    cost = layers.square_error_cost(current_Q1, target_Q) + layers.square_error_cost(current_Q2, target_Q)
    cost = layers.reduce_mean(cost)
    optimizer = fluid.optimizer.AdamOptimizer(self.critic_lr)
    optimizer.minimize(cost)
    return cost
```

buffer: 1e6 batch: 256 γ : 0.95 tau: 0.005 $A(s) \in (-\text{max}, \text{max})$
 actor-lr: $1e-3$ critic-lr: $1e-3$ fc1: 300 fc2: 300
 DU: noise: $\mu=0, \sigma=0.15, \sigma=0.2$
 就用这参数吧

(两次 Critic 更新, 行为策略才更新一次)
 For TD3: $a' = \text{clip}(V_{\phi}(s') + \text{noise})$ $V_1', V_2' = Q_{\phi}(s', a')$
 $y = r + \gamma \cdot \min V_1', V_2'$ $\tilde{y}_1, \tilde{y}_2 = Q_{\phi}(s, a)$
 $\text{cost} = \text{MSE}(\tilde{y}_1, y) + \text{MSE}(\tilde{y}_2, y)$
 如果在前 20 万 step 训练中, 性能与数据没有上到 1000, train 与数据没上到 0 就放弃
 行为策略: $\text{clip}(a + \eta(0, 0.1 a_{\max}), -a_{\max}, a_{\max})$
 在价值函数更新目标上, 也增加了扰动, 使价值函数更 robust

Critic 可以相比 actor 学慢一点, SAC 貌似最差, 可以和 DDPG 比一下,