Day 4. Soft Actor-Critic (SAC)

NPEX Reinforcement Learning

2021.07.29 Jaeuk Shin, Minkyu Park



SAC - Review



SAC - Review

How to incentivize exploration?

idea: augment reward as follows:

$$\sum_{t=0}^{T-1} \left(r(s_t, a_t) + \alpha \mathcal{H}(\pi(\cdot|s_t)) \right)$$

where $\mathcal{H}(\pi(\cdot|s_t))$: **entropy** of action distribution $\pi(\cdot|s_t)$

How to solve new MDP with this novel reward criterion?

 \rightarrow soft Q-learning, soft Bellman equation, etc.





Gaussian actor network $a \sim \mathcal{N}(\mu_{\phi}(s), \sigma_{\phi}(s))$

two critic networks $Q_1(s, a; \theta_1), Q_2(s, a; \theta_2)$

target & loss construction

Remind: Entropy
$$H(X) = -\sum_{x \in X} p(x) \log p(x) = E[-\log p(x)]$$

$$y_j = r_j + \gamma \mathbb{E}_{a \sim \pi(\cdot|s_j)} \left(Q(s'_j, a) - \alpha \log \pi(a|s_j) \right)$$

$$\to y_j = r_j + \gamma \mathbb{E}_{a \sim \pi_{\phi^-}(\cdot | s_j)} \left(\min_{i=1,2} Q_i(s'_j, a; \theta_i^-) - \alpha \log \pi_{\phi^-}(a | s_j) \right)$$



actor loss:

$$\alpha \log \pi_{\phi}(f_{\phi}(\epsilon_j, s_j)|s_j) - \min_{i=1,2} Q_i(s_j, f_{\phi}(\epsilon_j, s_j))$$

Remark. π_{ϕ} : probability density, f_{ϕ} : actor network



```
actor definition - 1<sup>st</sup> step
def forward(self, state, eval=False, with log prob=False):
    x = F.relu(self.fc1(state))
    x = F.relu(self.fc2(x))
    mu = self.fc3(x)
    log_sigma = self.fc4(x)
    # clip value of log_sigma, as was done in Haarnoja's implementation of SAC:
    # https://github.com/haarnoja/sac.git
    log_sigma = torch.clamp(log_sigma, -20.0, 2.0)
                          what we are doing here: compute distribution params \mu_{\phi}(s) and \sigma_{\phi}(s)
    sigma = torch.exp(log_sigma)
    distribution = Independent(Normal(mu, sigma), 1)
```

actor definition - 2nd step

```
if not eval:
   # use rsample() instead of sample(), as sample() does not allow back-propagation through params
    u = distribution.rsample()
   if with log prob:
                                                 Reparameterization trick to ease computation of \nabla \log \pi(a|s)
        log prob = distribution.log prob(u)
        \log_{prob} = 2.0 * torch.sum((np.log(2.0) + 0.5 * np.log(self.ctrl_range) - u - F.softplus(-2.0 * u)), dim=1)
                                                                                   softplus(x) = log(1 + e^x)
   else:
       log prob = None
else:
    u = mu
   log_prob = None
# apply tanh so that the resulting action lies in (-1, 1)^D
a = self.ctrl range * torch.tanh(u)
                                           u \sim \mathcal{N}(\mu_{\phi}(s), \sigma_{\phi}(s)) \longrightarrow a = \tanh(u) \sim ?
return a, log prob
```



```
def init (self, dimS, dimA, hidden1, hidden2):
    super(DoubleCritic, self).__init__()
    self.fc1 = nn.Linear(dimS + dimA, hidden1)
    self.fc2 = nn.Linear(hidden1, hidden2)
    self.fc3 = nn.Linear(hidden2, 1)
    self.fc4 = nn.Linear(dimS + dimA, hidden1)
    self.fc5 = nn.Linear(hidden1, hidden2)
    self.fc6 = nn.Linear(hidden2, 1)
def forward(self, state, action):
   x = torch.cat([state, action], dim=1)
   x1 = F.relu(self.fc1(x))
   x1 = F.relu(self.fc2(x1))
   x1 = self.fc3(x1)
   x2 = F.relu(self.fc4(x))
   x2 = F.relu(self.fc5(x2))
   x2 = self.fc6(x2)
```

```
def Q1(self, state, action):
    x = torch.cat([state, action], dim=1)
    x = F.relu(self.fc1(x))
    x = F.relu(self.fc2(x))
    x = self.fc3(x)
```

this is how we define twin critics!



```
with torch.no grad():
            # Get action with log(pi(a|s)) (also gradient)
            next_actions, log_probs = agent.pi(next_obs_batch, with_log_prob=True)
            # To calculate TQ, we need Q(s',pi(s'))
            target_q1, target_q2 = agent.target_Q(next_obs_batch, next_actions)
            # To mitigate overestimation! - Idea from TD3
                                                                 trick! (why?)
 8
            target_q = torch.min(target_q1, target_q2)
 9
10
            # TQ^pi = r + gamma [ Q(s',pi(s')) - alpha H(pi(s')) ]
11
            # Recall : H = sum[ -P(X) * log(P(x)) ] = E [ -log(P(x)) ]
12
            # Recall : H \approx -log(P(x))
13
            TQ = rew_batch + agent.gamma * masks * (target_q - agent.alpha * log probs)
14
15
         # Calculate MSELoss
16
17
        Q1, Q2 = agent.Q(obs_batch, act_batch)
        Q_{loss1} = torch.mean((Q1 - TQ)**2)
        Q loss2 = torch.mean((Q2 - TQ)**2)
19
        Q loss = Q loss1 + Q loss2
20
21
22
         # Gradient descent
         agent.Q optimizer.zero grad()
23
        Q loss.backward()
24
         agent.Q optimizer.step()
25
```

← training critics

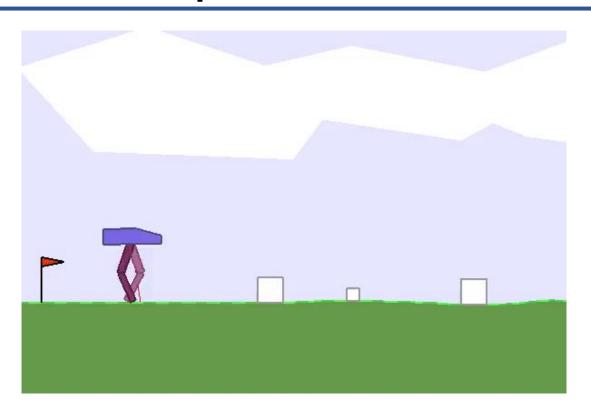


```
actions, log probs = agent.pi(obs batch, with log prob=True)
         freeze(agent.Q)
         q1, q2 = agent.Q(obs_batch, actions)
         q = torch.min(q1, q2)
 6
         # Need to perform gradient ascent, so (-) is required
                                                                                 training actor
         # TODO: build policy loss
 8
         #pi loss = torch.mean( #TODO )
 9
         pi loss = torch.mean(agent.alpha * log probs - q)
10
11
12
        # Gradient ascent
13
         agent.pi optimizer.zero grad()
         pi loss.backward()
14
         agent.pi_optimizer.step()
15
```

SAC - Experiment



SAC - Experiment



Is it successful on BipedalWalker-v3? (state dim: 24 / action dim: 4)



Thank you

