

PPO算法：

回顾TRPO

我们用老策略 π_{θ} 采样的数据来训练新策略 $\pi_{\theta'}$ ，所以新策略 $\pi_{\theta'}$ 距离不应该与老策略距离太远，如果差距过大，那么老策略采样出的数据会越来越没有价值，所以使用KL散度来约束策略更新的幅度

在经过约束后，我们可以近似认为前后两个策略的状态分布是一样的，但是在对动作上的选择不同的策略还是有区别的，因此使用**重要性采样**来让我们老策略采样出的历史的action，用来评估我们新的策略。从而去优化我们新的策略目标值

但TRPO在求解过程中存在一定的**不足**：

- 近似带来的误差

这也是重要性采样的通病，虽然我们使两个策略的KL散度保证不超过某个上界值，但即使如此我们会仍然碰到重要性采样估计不准，甚至过大的情况，导致学习的误差会比较大

- 求解约束优化问题的困难

另外一方面，就是近似的去求解约束优化问题本身的一些困难。TRPO 使用很多复杂的方法来进行求解，比如泰勒展开近似、共轭梯度、线性搜索等方法，这使计算过程非常复杂，每一步更新的运算量非常大。

在TRPO这些痛点的基础上进行改进——一种新的算法PPO算法

PPO全称**Proximal Policy Optimization**

PPO中的第一个P为Proximal，也就是近似替代的意思，它在TRPO上有3个改进

PPO—截断

第一项改进，就是针对**重要性采样**来进行改进，PPO提出了一种clip（截断;修建）的方式，将重要性采样的权重，ratio本身去做了一个截断，即下式子中的 $r_t(\theta)$ ：

$$L^{CPI}(\theta) = \hat{\mathbb{E}}_t \left[\frac{\pi_{\theta}(a_t|s_t)}{\pi_{\theta_{old}}(a_t|s_t)} \hat{A}_t \right] = \hat{\mathbb{E}}_t [r_t(\theta) \hat{A}_t]$$

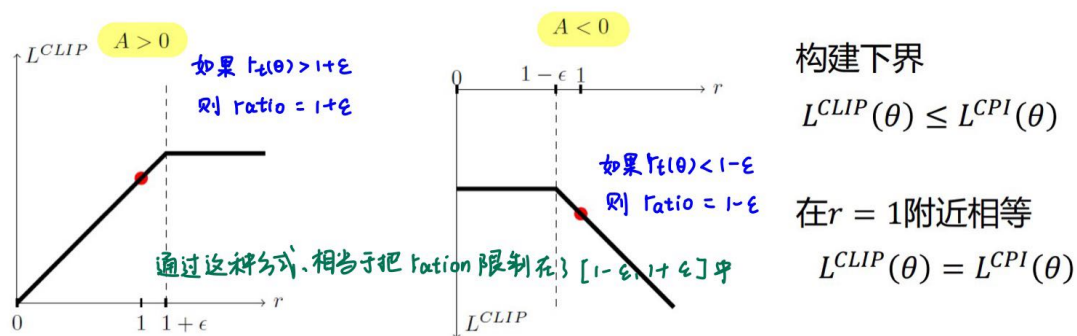
具体截断的目标值为：

$$L^{CLIP}(\theta) = \hat{\mathbb{E}}_t [\min(r_t(\theta) \hat{A}_t, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}_t)]$$

它是两项值，然后取最小值的关系。

- 第一个是原来值，TRPO中的ratio乘以优势函数
- 第二项，也就是该算法提出的基于截断式的优化

重点看第二项，第二项将原来的ratio，与 $1+\epsilon$ 和 $1-\epsilon$ 做一个上下界的截断



如果不进行截断, ratio会变得非常的不确定, 比如我们的ratio可能会存在大于 $1+\epsilon$ 的情况, 即老策略采样到的动作的概率作太小, 但是它作为分母就会导致ratio变得很大, 这样就会使得权重变得过大, 由于我们限制了两个策略前后的KL divergence, 所以这两个策略的ratio就不应该有过大或者过下的情况出现

所以通过截断的方式, 就使TRPO的鲁棒性或者说是稳定性获得了提升

PPO—优势函数 A_t

PPO在TRPO上的第二个改进, 即优势函数使用多步时序差分

$$\hat{A}_t = -V(s_t) + r_t + \gamma r_{t+1} + \dots + \gamma^{T-t+1} r_{T-1} + \gamma^{T-t} V(s_T)$$

在一般的优势函数中, 动作价值函数直接就是当前奖励加上一步的下一个状态价值乘以折扣然后与当前状态价值函数做差来作为优势函数, 这是一步的时序差分计算出来的优势函数

而PPO中, 往前看一定数量的步数, 使得我们获得了T步的经验数据作为整个优势函数的信息量, 显然是比一步带来的信息量会更大, 使得我们目标的估计会更加靠谱一些

- 在每次迭代中, 并行 N 个actor收集 T 步经验数据
- 计算每步的 \hat{A}_t 和 $L^{CLIP}(\theta)$, 构成mini-batch
- 更新参数 θ , 并更新 $\theta_{old} \leftarrow \theta$

PPO—惩罚

PPO相比TRPO上的第三个改进为KL散度惩罚项参数自适应的方法

在TRPO中的惩罚项系数为 λ , 如果策略的更新幅度增大, λ 也会增大, 从而增加惩罚力度, 来约束策略的变化幅度

λ 的更新是通过KL散度进行更新的, 如果这个KL散度超过了 ϵ , 即 $D_{KL} - \epsilon > 0$, 就加上这个超出的部分 $D_{KL} - \epsilon$ 然后乘以某个比率来更新这个 λ , 即 λ 的更新为:

$$\lambda \leftarrow \lambda + \alpha(D_{KL}(\pi_{\theta'}(a_t|s_t) \parallel \pi_{\theta}(a_t|s_t)) - \epsilon)$$

而这个加法型的惩罚项更新可能会比较慢, 以至于这个约束并不能够及时满足目标更新 (? 为啥加法就慢啊

所以PPO利用乘和除的方法, 来改进惩罚项参数使得PPO中的 β 值, 也就是TRPO中的 λ 值。它的改进速度会快很多, 具体的改进方法为:

- 计算KL值 $d = \hat{\mathbb{E}}_t [\text{KL}[\pi_{\theta_{\text{old}}}(\cdot | s_t) | \pi_{\theta}(\cdot | s_t)]]$ 计算每一时刻的KL divergence 具体的期望值d
 - a) 如果 $d < d_{\text{targ}}/1.5$, 更新 $\beta \leftarrow \beta/2$ 如果期望值d比目标值要小且小1.5倍以上, 那么我们认为此时处于一个比较安全的区域, 所以我们可以减小我们的惩罚力度, 使 $\beta/2 \rightarrow \beta$
 - b) 如果 $d > d_{\text{targ}} \times 1.5$, 更新 $\beta \leftarrow \beta \times 2$ 若d大于目标值的1.5倍, 说明我们Dkl过大, 所以需要加大惩罚力度, $2 \times \beta \rightarrow \beta$
- 事先设定的超参数 经验参数

通过这种方式, 使得我们对 β 的管控, 是否超过当前目标值 (用于限制学习策略和之前一轮策略的差距) 来使我们的 β 更加敏感, 其作用更加明显一些

这就是PPO相比TRPO的第三点改进

代码实践:

```

1  class PPO:
2      ''' PPO算法, 采用截断方式 '''
3      def __init__(self, state_dim, hidden_dim, action_dim, actor_lr,
4          critic_lr,
5              lambda, epochs, eps, gamma, device):
6          self.actor = PolicyNet(state_dim, hidden_dim, action_dim).to(device)
7          self.critic = ValueNet(state_dim, hidden_dim).to(device)
8          self.actor_optimizer = torch.optim.Adam(self.actor.parameters(),
9              lr=actor_lr)
10         self.critic_optimizer = torch.optim.Adam(self.critic.parameters(),
11             lr=critic_lr)
12         self.gamma = gamma
13         self.lambda = lambda
14         self.epochs = epochs # 一条序列的数据用来训练轮数
15         self.eps = eps # PPO中截断范围的参数
16         self.device = device
17
18     def take_action(self, state):
19         state = torch.tensor([state], dtype=torch.float).to(self.device)
20         probs = self.actor(state)
21         action_dist = torch.distributions.Categorical(probs)
22         action = action_dist.sample()
23         return action.item()
24
25     def update(self, transition_dict):
26         states = torch.tensor(transition_dict['states'],
27             dtype=torch.float).to(self.device)
28         actions = torch.tensor(transition_dict['actions']).view(-1, 1).to(
29             self.device)
30         rewards = torch.tensor(transition_dict['rewards'],
31             dtype=torch.float).view(-1,
32             1).to(self.device)
33         next_states = torch.tensor(transition_dict['next_states'],
34             dtype=torch.float).to(self.device)
35         dones = torch.tensor(transition_dict['dones'],
36             dtype=torch.float).view(-1, 1).to(self.device)
37         td_target = rewards + self.gamma * self.critic(next_states) * (1 -
38             dones)
39         td_delta = td_target - self.critic(states)
40         advantage = rl_utils.compute_advantage(self.gamma, self.lambda,

```

```

39 td_delta.cpu()).to(self.device)
40     old_log_probs = torch.log(self.actor(states).gather(1,
41
42     actions)).detach()
43
44     for _ in range(self.epochs):
45         log_probs = torch.log(self.actor(states).gather(1, actions))
46         ratio = torch.exp(log_probs - old_log_probs)
47         surr1 = ratio * advantage
48         surr2 = torch.clamp(ratio, 1 - self.eps,
49                             1 + self.eps) * advantage # 截断
50         actor_loss = torch.mean(-torch.min(surr1, surr2)) # PPO损失函数
51         critic_loss = torch.mean(
52             F.mse_loss(self.critic(states), td_target.detach()))
53         self.actor_optimizer.zero_grad()
54         self.critic_optimizer.zero_grad()
55         actor_loss.backward()
56         critic_loss.backward()
57         self.actor_optimizer.step()
58         self.critic_optimizer.step()

```

深度策略梯度算法小结

- 相比价值函数学习最小化TD误差的目标，策略梯度方法直接优化策略价值的目标 更加贴合强化学习本质目标
- 基于神经网络的策略在优化时容易因为一步走得太大而变得很差，进而下一轮产生 很低质量的经验数据，进一步无法学习好
- Trust Region一类方法限制一步更新前后策略的差距（用KL散度），进而对策略 价值做稳步地提升
- PPO在TRPO的基础上进一步通过限制importance ratio的范围，构建优化目标 的下界，进一步保证优化的稳定效果，是目前最常用的深度策略梯度算法
- 针对连续动作的决定性策略，可以从构建的critic中直接回传梯度到动作上，然后 通过链式法则进一步将梯度回传到策略网络中
- 分布式的actor-critic算法能够充分利用多核CPU资源采样环境的经验数据，利用 GPU资源异步地更新网络，这有效提升了DRL的训练效率