

# 彪哥带你学强化学习

## 15.深入理解TRPO算法(2)

DEEPLY UNDERSTAND REINFORCEMENT LEARNING

讲师：韩路彪





### 理论

- 每训练一步，整体回报都有提升
- 目标函数 $\eta$ ，近似函数 $L$
- 如果能保证  $|L_{\pi} - \eta| \leq X$ ，就能在  $\nabla L_{\pi} \geq X$  前提下做到  $\nabla \eta \geq 0$

### 理论实现

- 给出  $|L_{\pi} - \eta| \leq X$  里边的上限 $X$

### 工程实现

- 实现每训练一步 $\eta$ 都有提升



$$|L_{\pi}(\tilde{\pi}) - \eta(\tilde{\pi})| \leq \frac{4\alpha^2\gamma\epsilon}{(1-\gamma)^2}$$

$$\alpha = D_{TV}^{max}(\pi_{old}, \pi_{new})$$

$$D_{TV}(p||q) = \frac{1}{2} \sum_i |p_i - q_i|$$

$$D_{TV}^{max}(\pi||\tilde{\pi}) = \max_s D_{TV}(\pi(.|s), \tilde{\pi}(.|s))$$

$$\epsilon = \max_{s,a} |A_{\pi}(s, a)|$$

$$\eta(\tilde{\pi}) \geq L_{\pi}(\tilde{\pi}) - \frac{4\alpha^2\gamma\epsilon}{(1-\gamma)^2}$$



## TRPO算法 —— 理论实现

$$\text{证明 } |L_{\pi}(\tilde{\pi}) - \eta(\tilde{\pi})| \leq \frac{4\alpha^2\gamma\epsilon}{(1-\gamma)^2}$$

$$\begin{aligned} & |L_{\pi}(\tilde{\pi}) - \eta(\tilde{\pi})| \\ &= \sum_{t=0}^{\infty} \gamma^t |\mathbb{E}_{\tau \sim \tilde{\pi}}[\bar{A}(s_t)] - \mathbb{E}_{\tau \sim \pi}[\bar{A}(s_t)]| \\ &= \sum_{t=0}^{\infty} \gamma^t |P(n_t = 0) \mathbb{E}_{s_t \sim \tilde{\pi} | n_t=0}[\bar{A}(s_t)] + P(n_t > 0) \mathbb{E}_{s_t \sim \tilde{\pi} | n_t>0}[\bar{A}(s_t)] \\ &\quad - (P(n_t = 0) \mathbb{E}_{s_t \sim \pi | n_t=0}[\bar{A}(s_t)] + P(n_t > 0) \mathbb{E}_{s_t \sim \pi | n_t>0}[\bar{A}(s_t)])| \\ &= \sum_{t=0}^{\infty} \gamma^t |P(n_t > 0) \mathbb{E}_{s_t \sim \tilde{\pi} | n_t>0}[\bar{A}(s_t)] - P(n_t > 0) \mathbb{E}_{s_t \sim \pi | n_t>0}[\bar{A}(s_t)]| \\ &\leq \sum_{t=0}^{\infty} \gamma^t P(n_t > 0) (|\mathbb{E}_{s_t \sim \tilde{\pi} | n_t>0}[\bar{A}(s_t)]| + |\mathbb{E}_{s_t \sim \pi | n_t>0}[\bar{A}(s_t)]|) \end{aligned}$$

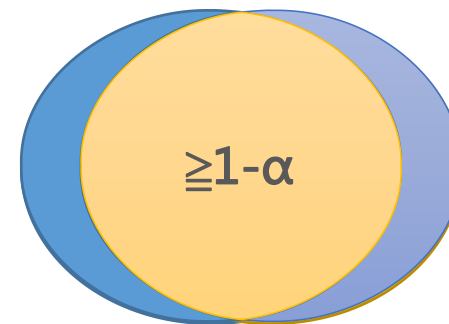
$n_t$ : 前 $n$ 步，有几步不一样



## TRPO算法 —— 理论实现

$$\begin{aligned}
 & |L_{\pi}(\tilde{\pi}) - \eta(\tilde{\pi})| \\
 & \leq \sum_{t=0}^{\infty} \gamma^t P(n_t > 0) (|\mathbb{E}_{s_t \sim \tilde{\pi} | n_t > 0} [\bar{A}(s_t)]| + |\mathbb{E}_{s_t \sim \pi | n_t > 0} [\bar{A}(s_t)]|) \\
 & \leq \sum_{t=0}^{\infty} \gamma^t 2(1 - (1 - \alpha)^t) 2\alpha \max_{s,a} |A(s, a)| \\
 & = 4\epsilon\alpha \sum_{t=0}^{\infty} \gamma^t (1 - (1 - \alpha)^t) \\
 & = 4\epsilon\alpha \left( \frac{1}{1 - \gamma} - \frac{1}{1 - \gamma(1 - \alpha)} \right) \\
 & = \frac{4\epsilon\alpha^2\gamma}{(1 - \gamma)(1 - \gamma(1 - \alpha))} \\
 & \leq \frac{4\epsilon\alpha^2\gamma}{(1 - \gamma)^2}
 \end{aligned}$$

对任意  $(a, \tilde{a}) | s \sim (\pi, \tilde{\pi})$  有  $P(a \neq \tilde{a} | s) \leq \alpha$



$$\epsilon = \max_{s,a} |A(s, a)|$$

等比级数