

# TRPO知识总结

- 随机策略梯度方法的问题：不好选取合适的步长；若步长太长，策略很容易发散；若步长太短，策略收敛速度很慢
- TRPO解决的问题：选取合适的步长，使新的策略的回报函数的值比旧策略大
- 定义策略的回报函数

$$\eta(\tilde{\pi}) = E_{\tau|\tilde{\pi}} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t) \right]$$

◦  $\tau: s_0, u_0, \dots, s_H, u_H$

- 将回报函数写成老策略的回报+某一大于零的项

$$\eta(\tilde{\pi}) = \eta(\pi) + E_{\tau|\tilde{\pi}} \left[ \sum_{t=0}^{\infty} \gamma^t (A_{\pi}(s_t, a_t)) \right]$$

◦ 证明过程

The image shows a handwritten derivation of the advantage function  $A_{\pi}(s_t, a_t)$ . The steps are as follows:

$$\begin{aligned} & E_{\tau|\tilde{\pi}} \left[ \sum_{t=0}^{\infty} \gamma^t A_{\pi}(s_t, a_t) \right] \\ &= E_{\tau|\tilde{\pi}} \left[ \sum_{t=0}^{\infty} \gamma^t (r(s_t) + \gamma V^{\tilde{\pi}}(s_{t+1}) - V^{\tilde{\pi}}(s_t)) \right] \\ &= E_{\tau|\tilde{\pi}} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t) + \sum_{t=0}^{\infty} \gamma^t (\gamma V^{\tilde{\pi}}(s_{t+1}) - V^{\tilde{\pi}}(s_t)) \right] \\ &= E_{\tau|\tilde{\pi}} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t) + (-V^{\tilde{\pi}}(s_0)) \right] \\ &= E_{\tau|\tilde{\pi}} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t) \right] - E_{s_0} [-V^{\tilde{\pi}}(s_0)] \\ &= \eta(\tilde{\pi}) - V(s_0) = \eta(\tilde{\pi}) - \eta(\pi) \end{aligned}$$

- 将优势函数的期望展开，写成对分别对状态空间和动作空间的积分

$$\eta(\tilde{\pi}) = \eta(\pi) + \sum_s \rho_{\tilde{\pi}}(s) \sum_a \tilde{\pi}(a|s) A_{\pi}(s, a)$$

- TRPO技巧1：忽略状态分布的变化，消除状态分布对新策略的依赖 - 对代价函数的近似1

$$L_{\pi}(\tilde{\pi}) = \eta(\pi) + \sum_s \rho_{\pi}(s) \sum_a \tilde{\pi}(a|s) A_{\pi}(s, a)$$

- TRPO技巧2：利用重要性采样处理动作分布，克服新的策略（动作分布）未知的问题  
采样策略：旧策略

$$L_{\pi}(\tilde{\pi}) = \eta(\pi) + E_{(s,a) \sim \pi} \left[ \frac{\tilde{\pi}(a|s)}{\pi(a|s)} A_{\pi}(s, a) \right]$$

- 替代回报函数和原回报函数在策略 $\pi_{old}$ 处一阶近似的证明

TRPO②

$$\begin{aligned}
 L_{\pi_{old}}(\pi_{old}) &= \eta(\pi_{old}) + E_{(s,a) \sim \pi_{old}} [A_{\pi_{old}}(s,a)] \\
 &= \eta(\pi_{old}) + \sum_s p_{\pi_{old}}(s) \sum_a \pi_{old}(a|s) \cdot A_{\pi_{old}}(s,a) \\
 &= \eta(\pi_{old}) + \sum_s p_{\pi_{old}}(s) \cdot \left( \sum_a \pi_{old}(a|s) Q_{\pi_{old}}(s,a) - \sum_a \pi_{old}(a|s) V_{\pi_{old}}(s) \right) \\
 &= \eta(\pi_{old}) + \sum_s p_{\pi_{old}}(s) \cdot (V_{\pi_{old}}(s) - V_{\pi_{old}}(s)) \\
 &= \eta(\pi_{old})
 \end{aligned}$$

$$\begin{aligned}
 \nabla_{\theta} L_{\pi_{old}}(\pi_{old}) \Big|_{\theta=\theta_{old}} &= \nabla_{\theta} \eta(\pi_{old}) + \nabla_{\theta} \left[ \sum_s p_{\pi_{old}}(s) \sum_a \pi_{old}(a|s) A_{\pi_{old}}(s,a) \right] \Big|_{\theta=\theta_{old}} \\
 &= 0 + \nabla_{\theta} \sum_s p_{\pi_{old}}(s) \sum_a \pi_{old}(a|s) \cdot (Q_{\pi_{old}}(s,a) - V_{\pi_{old}}(s)) \\
 &= \sum_s p_{\pi_{old}}(s) (\nabla_{\theta} \sum_a \pi_{old}(a|s) Q_{\pi_{old}}(s,a) - \nabla_{\theta} V_{\pi_{old}}(s)) \\
 &= \sum_s p_{\pi_{old}}(s) \sum_a \nabla_{\theta} \pi_{old}(a|s) Q_{\pi_{old}}(s,a) \Big|_{\theta=\theta_{old}}
 \end{aligned}$$

$$\begin{aligned}
 \nabla_{\theta} \eta(\pi_{old}) \Big|_{\theta=\theta_{old}} &= \nabla_{\theta} E_{(s,a) \sim \pi_{old}} \left[ \sum_{t=0}^{\infty} \gamma^t R_{t+1} \right] = \nabla_{\theta} V_{\pi_{old}}(s_0) \Big|_{\theta=\theta_{old}} \\
 \frac{\partial V_{\pi_{old}}(s)}{\partial \theta} &= \frac{\partial}{\partial \theta} \sum_a \pi_{old}(a|s) Q_{\pi_{old}}(s,a) = \sum_a \left[ \frac{\partial \pi_{old}(a|s)}{\partial \theta} Q_{\pi_{old}}(s,a) + \pi_{old}(a|s) \frac{\partial Q_{\pi_{old}}(s,a)}{\partial \theta} \right] \\
 &\equiv \sum_a \left[ \frac{\partial \pi_{old}(a|s)}{\partial \theta} Q_{\pi_{old}}(s,a) + \pi_{old}(a|s) \sum_{s'} \gamma p_{ss'}^a \frac{\partial V_{\pi_{old}}(s')}{\partial \theta} \right] \\
 &= \sum_a \sum_{s'} \gamma p_{ss'}^a \pi_{old}(a|s) \frac{\partial \pi_{old}(s'|s)}{\partial \theta} Q_{\pi_{old}}(s',a)
 \end{aligned}$$

$$\nabla_{\theta} \eta(\pi_{old}) \Big|_{\theta=\theta_{old}} = \sum_s p_{\pi_{old}}(s) \sum_a \frac{\partial \pi_{old}(a|s)}{\partial \theta} Q_{\pi_{old}}(s,a) \Big|_{\theta=\theta_{old}} = \sum_s p_{\pi_{old}}(s) \sum_a \nabla_{\theta} \pi_{old}(a|s) \cdot Q_{\pi_{old}}(s,a)$$

- 上图中红色一行推导的理解：上一行表示对于 $V_{\pi_{old}}(s)$ 的梯度，可以转化成对策略梯度 \*  $Q_{\pi_{old}}(s,a)$  + 所有可能的下一状态的 $V_{\pi}(s')$ 的梯度 \* 这一状态 $s'$ 出现的概率，因此是一种迭代的计算方法，继续迭代下去，会遍历到所有可能的状态与行为；因此，在红色行中，计算策略梯度 \* 行为值函数的期望：首先对行为空间进行积分，之后对状态 $s$ 进行积分；对状态 $s$ 积分时，概率是从 $s$ 经过任意步数转变为 $x$ 的概率的和，因为是要计算准确的值函数梯度，必须要计算到结束。
- 单调性的证明：书上不等式 (8.8)
- 为决定更新步长，问题转化为

$$\begin{aligned}
 &\underset{\theta}{\text{maximize}} E_{s \sim \rho_{\theta_{old}}, a \sim \pi_{\theta_{old}}} \left[ \frac{\pi_{\theta}(a|s)}{\pi_{\theta_{old}}(a|s)} A_{\theta_{old}}(s,a) \right] \\
 &\text{subject to } D_{KL}^{\max}(\theta_{old}, \theta) \leq \delta
 \end{aligned}$$

- TRPO第三个技巧：利用平均KL散度代替最大KL散度
  - 最大：遍历所有 $s$ ，对每个 $s$ 都有一个KL散度，找最大的那个
  - 平均： $s$ 分布服从 $\rho$ ，对 $s$ 的分布求期望，概率是 $\rho(s)$ ，值是这个 $s$ 下的KL散度
- TRPO第四个技巧： $s \sim \rho_{\theta_{old}} \rightarrow s \sim \pi_{\theta_{old}}$ 
  - 理解：忽略折扣因子？
- 因此，平均KL散度的期望中， $s$ 服从的分布也变为 $\pi_{\theta_{old}}$ 
  - 约束条件变为

$$E_{s \sim \pi_{old}} [D_{KL}(\pi_{old} || \pi)]$$

- 采样，利用样本均值代替期望 - ppt44

- 对TRPO目标函数一阶逼近、对约束条件二阶逼近

目标函数一阶泰勒展开  $L(\theta) = E_{s \sim \pi_{old}, a \sim \pi_{old}} \left[ \frac{\pi_{old}(a|s)}{\pi_{old}(a|s)} A_{\theta_{old}}(s, a) \right]$

~~在  $\theta_{old}$  处~~

$$L(\theta) = \frac{1}{0!} L(\theta_{old}) + \frac{1}{1!} \nabla_{\theta} L(\theta) \cdot (\theta - \theta_{old})$$

$$= 1 \cdot E_{s, a \sim \pi_{old}} [A_{\theta_{old}}(s, a)] + \nabla_{\theta} L(\theta) \Big|_{\theta = \theta_{old}} \cdot (\theta - \theta_{old})$$

$$= \nabla_{\theta} L(\theta) \Big|_{\theta = \theta_{old}} (\theta - \theta_{old})$$

约束条件在  $\theta_{old}$  处二阶泰勒展开

$$D_{KL}[\pi_{\theta_{old}} || \pi_{\theta}] \approx E_{s \sim \pi_{old}} \log \pi_{\theta_{old}} - E_{s \sim \pi_{\theta}} \log \pi_{\theta}$$

$$= E_{\theta_{old}} [\log \pi_{\theta_{old}}] - (E_{\theta_{old}} [\log \pi_{\theta_{old}}] + \frac{1}{1!} E_{\theta_{old}} [\nabla_{\theta} \log \pi_{\theta}] \cdot \Delta \theta + \frac{1}{2!} E_{\theta_{old}} [\nabla_{\theta}^2 \log \pi_{\theta}] \cdot (\theta - \theta_{old})^2)$$

$$= -\frac{1}{2} E_{\theta_{old}} [\nabla_{\theta}^2 \log \pi_{\theta}] \Big|_{\theta = \theta_{old}} \cdot (\theta - \theta_{old})^2$$

- KL散度非负（非负性证明见教材P155），在  $\theta = \theta_{old}$  处，KL散度=0，所以  $\theta_{old}$  是KL散度函数( $\theta$ )的一个极值点，因此KL散度函数一阶导数=0