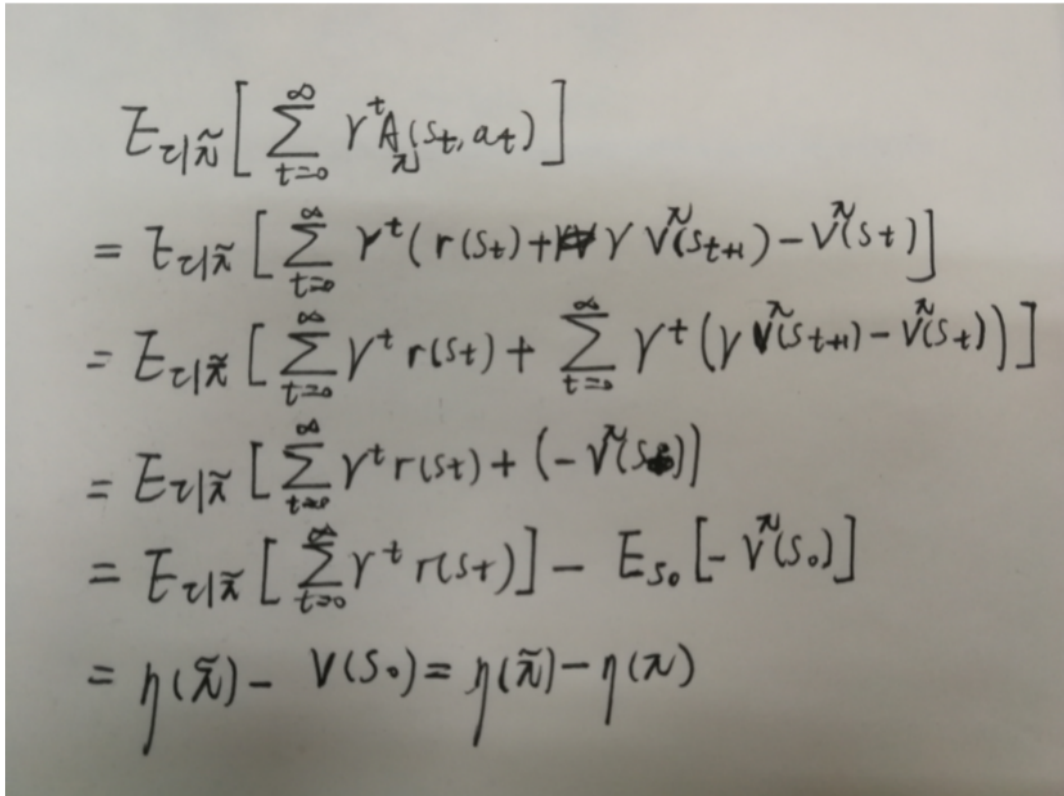


# 强化学习作业/第8次/个人作业/TRPO三个等式的证明

## 1 TRPO起点等式的证明

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

◦ 证明过程



Handwritten proof of the TRPO starting equation:

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

2 替代回报函数和原回报函数在 $\theta = \theta_{old}$ 处相等、一阶近似

- 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

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

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

$$\begin{aligned}
 \nabla_{\theta} \eta(\pi_{old})|_{\theta=\theta_{old}} &= \nabla_{\theta} E_{(s,a) \sim \pi_{old}} [r + \gamma V_{\pi_{old}}(s)] = \nabla_{\theta} V_{\pi_{old}}(s) \Big|_{\theta=\theta_{old}} \\
 \frac{\partial V_{\pi_{old}}(s)}{\partial \theta} &= \frac{\partial}{\partial \theta} \sum_a \pi_{old}(a|s) \cdot Q_{\pi_{old}}(s,a) = \sum_a \left[ \frac{\partial \pi_{old}(a|s)}{\partial \theta} \cdot Q_{\pi_{old}}(s,a) + \pi_{old}(a|s) \cdot \frac{\partial Q_{\pi_{old}}(s,a)}{\partial \theta} \right] \\
 &= \sum_a \left[ \frac{\partial \pi_{old}(a|s)}{\partial \theta} \cdot Q_{\pi_{old}}(s,a) + \pi_{old}(a|s) \cdot \sum_{s'} \gamma P_{ss'}^a \cdot \frac{\partial V_{\pi_{old}}(s')}{\partial \theta} \right] \\
 &= \sum_a \sum_{s'} \gamma P_{ss'}^a \cdot \pi_{old}(a|s) \cdot \frac{\partial \pi_{old}(a|s')}{\partial \theta} \cdot Q_{\pi_{old}}(s,a) \\
 \nabla_{\theta} \eta(\pi_{old})|_{\theta=\theta_{old}} &= \sum_s p_{\pi_{old}}(s) \sum_a \frac{\partial \pi_{old}(a|s)}{\partial \theta} \cdot Q_{\pi_{old}}(s,a) \Big|_{\theta=\theta_{old}} = \sum_s p_{\pi_{old}}(s) \sum_a \nabla_{\theta} \pi_{old}(a|s) \cdot Q_{\pi_{old}}(s,a)
 \end{aligned}$$

- 上图中红色一行推导的理解：上一行表示对于 $V_{\pi_0}(s)$ 的梯度，可以转化成对策略梯度 $* Q_{\pi_0}(s,a) +$ 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。

### 3 TRPO优化目标的一阶近似、约束条件的二阶近似

- 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  $L(\theta) = E_{(s,a) \sim \pi_{old}, a \sim \pi_{old}} \left[ \frac{\pi_{old}(a|s)}{\pi_{old}(a|s)} A_{\pi_{old}}(s,a) \right]$

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

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

约束条件在 $\theta_{old}$ 处二阶泰勒展开

$$\begin{aligned}
 D_{KL}[\pi_{old} || \pi_{\theta}] &= E_{s \sim \pi_{old}} \log \pi_{old} - E_{s \sim \pi_{\theta}} \log \pi_{\theta} \\
 &= E_{\pi_{old}} [\log \pi_{old}] - (E_{\pi_{old}} [\log \pi_{old}] + \frac{1}{1!} E_{\pi_{old}} [\nabla_{\theta} \log \pi_{\theta}] \cdot \Delta \theta \\
 &\quad + \frac{1}{2!} E_{\pi_{old}} [\nabla_{\theta}^2 \log \pi_{\theta}] \cdot (\theta - \theta_{old})^2) \\
 &= -\frac{1}{2} E_{\pi_{old}} [\nabla_{\theta}^2 \log \pi_{\theta}] \Big|_{\theta=\theta_{old}} (\theta - \theta_{old})^2
 \end{aligned}$$

- 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