

第九章 策略学习高级技巧

本章介绍策略学习的高级技巧。第 9.1 介绍置信域策略优化 (TRPO)，它是一种策略学习方法，可以代替策略梯度方法。第 9.2 介绍熵正则，可以用在所有的策略学习方法中。

9.1 Trust Region Policy Optimization (TRPO)

置信域策略优化 (Trust Region Policy Optimization, TRPO) 是一种策略学习方法，跟以前学的策略梯度有很多相似之处。跟策略梯度方法相比，TRPO 有两个优势：第一，TRPO 表现更稳定，收敛曲线不会剧烈波动，而且对学习率不敏感；第二，TRPO 用更少的经验（即智能体收集到的状态、动作、奖励）就能达到与策略梯度方法相同的表现。

学习 TRPO 的关键在于理解置信域方法 (Trust Region Methods)。置信域方法不是 TRPO 的论文提出的，而是数值最优化领域中一类经典的算法，历史至少可以追溯到 1970 年。TRPO 论文的贡献在于巧妙地把置信域方法应用到强化学习中，取得非常好的效果。

本节分以下 4 小节讲解 TRPO：第 9.1.1 小节介绍置信域方法，第 9.1.2 节回顾策略学习，第 9.1.3 节推导 TRPO，第 9.1.4 讲解 TRPO 的算法流程。

9.1.1 置信域方法

有这样一个优化问题： $\max_{\theta} J(\theta)$ 。这里的 $J(\theta)$ 是目标函数， θ 是优化变量。求解这个优化问题的目的是找到一个变量 θ 使得目标函数 $J(\theta)$ 取得最大值。有各种各样的优化算法用于解决这个问题。几乎所有的数值优化算法都是做这样的迭代：

$$\theta_{\text{new}} \leftarrow \text{Update}(\text{Data}; \theta_{\text{now}}).$$

此处的 θ_{now} 和 θ_{new} 分别是优化变量当前的值和新的值。不同算法的区别在于具体怎么样利用数据更新优化变量。

置信域方法用到一个概念——**置信域**。下面介绍置信域。给定变量当前的值 θ_{now} ，用 $\mathcal{N}(\theta_{\text{now}})$ 表示 θ_{now} 的一个邻域。举个例子：

$$\mathcal{N}(\theta_{\text{now}}) = \left\{ \theta \mid \|\theta - \theta_{\text{now}}\|_2 \leq \Delta \right\}. \quad (9.1)$$

这个例子中，集合 $\mathcal{N}(\theta_{\text{now}})$ 是以 θ_{now} 为球心、以 Δ 为半径的球；见右图。球中的点都足够接近 θ_{now} 。

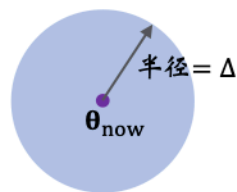


图 9.1: 公式(9.1)中的邻域 $\mathcal{N}(\theta_{\text{now}})$ 。

置信域方法需要构造一个函数 $L(\theta \mid \theta_{\text{now}})$ ，这个函数要满足这个条件：

$$L(\theta \mid \theta_{\text{now}}) \text{ 很接近 } J(\theta), \quad \forall \theta \in \mathcal{N}(\theta_{\text{now}}),$$

那么集合 $\mathcal{N}(\theta_{\text{now}})$ 就被称作**置信域**。顾名思义，在 θ_{now} 的邻域上，我们可以信任 $L(\theta \mid \theta_{\text{now}})$ ，可以拿 $L(\theta \mid \theta_{\text{now}})$ 来替代目标函数 $J(\theta)$ 。

图 9.2 用一个一元函数的例子解释 $J(\theta)$ 和 $L(\theta | \theta_{\text{now}})$ 的关系。图中横轴是优化变量 θ ，纵轴是函数值。如图 9.2(a) 所示，函数 $L(\theta | \theta_{\text{now}})$ 未必在整个定义域上都接近 $J(\theta)$ ，而只是在 θ_{now} 的邻域里接近 $J(\theta)$ 。 θ_{now} 的邻域就叫做置信域。

通常来说， J 是个很复杂的函数，我们甚至可能不知道 J 的解析表达式（比如 J 是某个函数的期望）。而我们人为构造出的函数 L 相对较为简单，比如 L 是 J 的蒙特卡洛近似，或者是 J 在 θ_{now} 这个点的二阶泰勒展开。既然可以信任 L ，那么不妨用 L 代替复杂的函数 J ，然后对 L 做最大化。这样比直接优化 J 要容易得多。这就是置信域方法的思想。具体来说，置信域方法做下面这两个步骤，一直重复下去，当无法让 J 的值增大的时候终止算法。

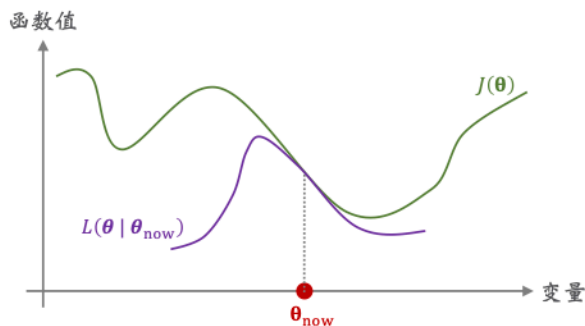
第一步——做近似：给定 θ_{now} ，构造函数 $L(\theta | \theta_{\text{now}})$ ，使得对于所有的 $\theta \in \mathcal{N}(\theta_{\text{now}})$ ，函数值 $L(\theta | \theta_{\text{now}})$ 与 $J(\theta)$ 足够接近。图 9.2(b) 解释了做近似这一步。

第二步——最大化：在置信域 $\mathcal{N}(\theta_{\text{now}})$ 中寻找变量 θ 的值，使得函数 L 的值最大化。把找到的值记作

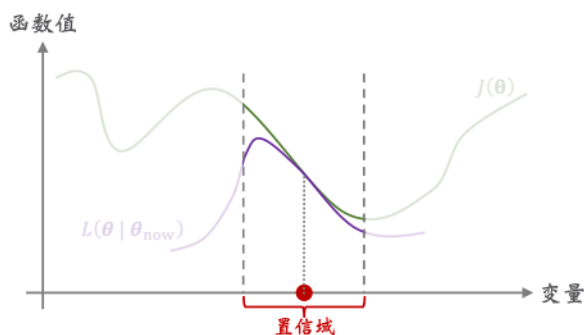
$$\theta_{\text{new}} = \operatorname{argmax}_{\theta \in \mathcal{N}(\theta_{\text{now}})} L(\theta | \theta_{\text{now}}).$$

图 9.2(c) 解释了最大化这一步。

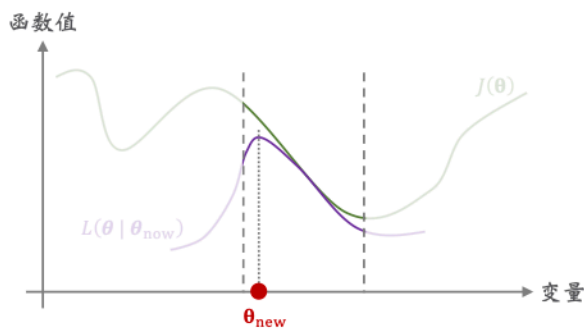
置信域方法其实是一类算法框架，而非一个具体的算法。有很多种方式实现置信域方法。第一步需要做近似，而做近似的方法多种多样，比如蒙特卡洛、二阶泰勒展开。第二步需要解一个带约束的最大化问题；求解这个问题又需要单独的数值优化算法，比如梯度投影算法、拉格朗日法。除此之外，置信域 $\mathcal{N}(\theta_{\text{now}})$ 也有多种多样的选择，既可以是球，也可以是两个概率分布的 KL 散度 (KL Divergence)，稍后会介绍。



(a) 构造 $L(\theta | \theta_{\text{now}})$ 作为 $J(\theta)$ 在点 θ_{now} 附近的近似。



(b) L 在点 θ_{now} 的邻域内接近 J ；这个领域就叫置信域。



(c) 在置信域内寻找最大化 L 的解，记作 θ_{new} 。

图 9.2: 一元函数的例子解释置信域和置信域算法。

9.1.2 策略学习

首先复习策略学习的基础知识。策略网络记作 $\pi(a|s; \theta)$ ，它是个概率密度函数。动作价值函数记作 $Q_\pi(s, a)$ ，它是回报的条件期望。状态价值函数记作

$$V_\pi(s) = \mathbb{E}_{A \sim \pi(\cdot|s; \theta)}[Q_\pi(s, A)] = \sum_{a \in \mathcal{A}} \pi(a|s; \theta) \cdot Q_\pi(s, a). \quad (9.2)$$

注意， $V_\pi(s)$ 依赖于策略网络 π ，所以依赖于 π 的参数 θ 。策略学习的目标函数是

$$J(\theta) = \mathbb{E}_S[V_\pi(S)]. \quad (9.3)$$

$J(\theta)$ 只依赖于 θ ，不依赖于状态 S 和动作 A 。第 7 章介绍的策略梯度方法（包括 REINFORCE 和 Actor-Critic）用蒙特卡洛近似梯度 $\nabla_\theta J(\theta)$ ，得到随机梯度，然后做随机梯度上升更新 θ ，使得目标函数 $J(\theta)$ 增大。

下面我们要把目标函数 $J(\theta)$ 变换成一种等价形式。从等式(9.2)出发，把状态价值写成

$$\begin{aligned} V_\pi(s) &= \sum_{a \in \mathcal{A}} \pi(a|s; \theta_{\text{now}}) \cdot \frac{\pi(a|s; \theta)}{\pi(a|s; \theta_{\text{now}})} \cdot Q_\pi(s, a) \\ &= \mathbb{E}_{A \sim \pi(\cdot|s; \theta_{\text{now}})} \left[\frac{\pi(A|s; \theta)}{\pi(A|s; \theta_{\text{now}})} \cdot Q_\pi(s, A) \right]. \end{aligned} \quad (9.4)$$

第一个等式很显然，因为连加中的第一项可以消掉第二项的分母。第二个等式把策略网络 $\pi(A|s; \theta_{\text{now}})$ 看做动作 A 的概率密度函数，所以可以把连加写成期望。由公式 (9.3) 与 (9.4) 可得定理 9.1。定理 9.1 是 TRPO 的关键所在，甚至可以说 TRPO 就是从这个公式推出的。

定理 9.1. 目标函数的等价形式

目标函数 $J(\theta)$ 可以等价写成：

$$J(\theta) = \mathbb{E}_S \left[\mathbb{E}_{A \sim \pi(\cdot|S; \theta_{\text{now}})} \left[\frac{\pi(A|S; \theta)}{\pi(A|S; \theta_{\text{now}})} \cdot Q_\pi(S, A) \right] \right].$$



公式中的期望是关于状态 S 和动作 A 求的。状态 S 的概率密度函数只有环境知道，而我们并不知道，但是我们可以从环境中获取 S 的观测值。动作 A 的概率密度函数是策略网络 $\pi(A|S; \theta_{\text{now}})$ ；注意，策略网络的参数是旧的值 θ_{now} 。

9.1.3 TRPO 数学推导

前面介绍了数值优化的基础和价值学习的基础，终于可以开始推导 TRPO。TRPO 是置信域方法在策略学习中的应用，所以 TRPO 也遵循置信域方法的框架，重复做近似和最大化这两个步骤，直到算法收敛。收敛指的是无法增大目标函数 $J(\theta)$ ，即无法增大期望回报。

第一步——做近似：我们从定理 9.1 出发。定理把目标函数 $J(\theta)$ 写成了期望的形式。我们无法直接算出期望，无法得到 $J(\theta)$ 的解析表达式；原因在于只有环境知道状态 S 的概率密度函数，而我们不知道。我们可以对期望做蒙特卡洛近似，从而把函数 J 近似成

函数 L 。用策略网络 $\pi(A|S; \theta_{\text{now}})$ 控制智能体跟环境交互，从头到尾玩完一局游戏，观测到一条轨迹：

$$s_1, a_1, r_1, s_2, a_2, r_2, \dots, s_n, a_n, r_n.$$

其中的状态 $\{s_t\}_{t=1}^n$ 都是从环境中观测到的，其中的动作 $\{a_t\}_{t=1}^n$ 都是根据策略网络 $\pi(\cdot|s_t; \theta_{\text{now}})$ 抽取的样本。所以，

$$\frac{\pi(a_t|s_t; \theta)}{\pi(a_t|s_t; \theta_{\text{now}})} \cdot Q_{\pi}(s_t, a_t) \quad (9.5)$$

是对定理 9.1 中期望的无偏估计。我们观测到了 n 组状态和动作，于是应该对公式 (9.5) 求平均，把得到均值记作：

$$L(\theta|\theta_{\text{now}}) = \frac{1}{n} \sum_{t=1}^n \underbrace{\frac{\pi(a_t|s_t; \theta)}{\pi(a_t|s_t; \theta_{\text{now}})} \cdot Q_{\pi}(s_t, a_t)}_{\text{定理 9.1 中期望的无偏估计}}. \quad (9.6)$$

既然连加里每一项都是期望的无偏估计，那么 n 项的均值 L 也是无偏估计。所以可以拿 L 作为目标函数 J 的蒙特卡洛近似。

公式 (9.6) 中的 $L(\theta|\theta_{\text{now}})$ 是对目标函数 $J(\theta)$ 的近似。可惜我们还无法直接对 L 求最大化，原因是我们不知道动作价值 $Q_{\pi}(s_t, a_t)$ 。解决方法是把 $Q_{\pi}(s_t, a_t)$ 近似成观测到的折扣回报：

$$u_t = r_t + \gamma \cdot r_{t+1} + \gamma^2 \cdot r_{t+2} + \dots + \gamma^{n-t} \cdot r_n.$$

拿 u_t 替代 $Q_{\pi}(s_t, a_t)$ ，那么公式 (9.6) 中的 $L(\theta|\theta_{\text{now}})$ 变成了

$$\tilde{L}(\theta|\theta_{\text{now}}) = \sum_{t=1}^n \frac{\pi(a_t|s_t; \theta)}{\pi(a_t|s_t; \theta_{\text{now}})} \cdot u_t. \quad (9.7)$$

总结一下，我们把目标函数 J 近似成 L ，然后又把 L 近似成 \tilde{L} 。

第二步——最大化： TRPO 把公式 (9.7) 中的 $\tilde{L}(\theta|\theta_{\text{now}})$ 作为对目标函数 $J(\theta)$ 的近似，然后求解这个带约束的最大化问题：

$$\max_{\theta} \tilde{L}(\theta|\theta_{\text{now}}); \quad \text{s.t. } \theta \in \mathcal{N}(\theta_{\text{now}}). \quad (9.8)$$

公式中的 $\mathcal{N}(\theta_{\text{now}})$ 是置信域，即 θ_{now} 的一个邻域。该用什么样的置信域呢？

- 一种方法是用以 θ_{now} 为球心、以 Δ 为半径的球作为置信域。这样的话，公式 (9.8) 就变成

$$\max_{\theta} \tilde{L}(\theta|\theta_{\text{now}}); \quad \text{s.t. } \|\theta - \theta_{\text{now}}\|_2 \leq \Delta. \quad (9.9)$$

- 另一种方法是用 KL 散度衡量两个概率密度函数—— $\pi(\cdot|s_i; \theta_{\text{now}})$ 和 $\pi(\cdot|s_i; \theta)$ ——的距离。两个概率密度函数区别越大，它们的 KL 散度就越大。反之，如果 θ 很接近 θ_{now} ，那么两个概率密度函数就越接近。用 KL 散度的话，公式 (9.8) 就变成

$$\max_{\theta} \tilde{L}(\theta|\theta_{\text{now}}); \quad \text{s.t. } \frac{1}{t} \sum_{i=1}^t \text{KL}[\pi(\cdot|s_i; \theta_{\text{now}}) \parallel \pi(\cdot|s_i; \theta)] \leq \Delta. \quad (9.10)$$

用球作为置信域的好处是置信域是简单的形状，求解最大化问题比较容易，但是用球做

置信域的实际效果不如用 KL 散度。

TRPO 的第二步——最大化——需要求解带约束的最大化问题 (9.9) 或者 (9.10)。注意，这种问题的求解并不容易；简单的梯度上升算法并不能解带约束的最大化问题。数值优化教材里面通常会有章节介绍这类问题的求解，有兴趣的话自己去阅读数值优化教材。这里就不详细解释如何求解问题 (9.9) 或者 (9.10)。读者可以这样看待优化问题：只要你能把一个优化问题的目标函数和约束条件写出来，通常会有数值算法能解决这个问题。

9.1.4 训练流程

在本节的最后，我们总结一下用 TRPO 训练策略网络的流程。TRPO 需要重复做近似和最大化这两个步骤：

1. 做近似——构造函数 \tilde{L} 近似目标函数 $J(\theta)$ ：

- (a). 设当前策略网络参数是 θ_{now} 。用策略网络 $\pi(a|s; \theta_{\text{now}})$ 控制智能体与环境交互，玩完一局游戏，记录下轨迹：

$$s_1, a_1, r_1, s_2, a_2, r_2, \dots, s_n, a_n, r_n.$$

- (b). 对于所有的 $t = 1, \dots, n$ ，计算折扣回报 $u_t = \sum_{k=t}^n \gamma^{k-t} \cdot r_k$ 。

- (c). 得出近似函数：

$$\tilde{L}(\theta | \theta_{\text{now}}) = \sum_{t=1}^n \frac{\pi(a_t | s_t; \theta)}{\pi(a_t | s_t; \theta_{\text{now}})} \cdot u_t.$$

2. 最大化——用某种数值算法求解带约束的最大化问题：

$$\theta_{\text{new}} = \underset{\theta}{\operatorname{argmax}} \tilde{L}(\theta | \theta_{\text{now}}); \quad \text{s.t. } \|\theta - \theta_{\text{now}}\|_2 \leq \Delta.$$

此处的约束条件是二范数距离。可以把它替换成 KL 散度，即公式 (9.10)。

TRPO 中有两个需要调的超参数：一个是置信域的半径 Δ ，另一个是求解最大化问题的数值算法的学习率。通常来说， Δ 在算法的运行过程中要逐渐缩小。虽然 TRPO 需要调参，但是 TRPO 对超参数的设置并不敏感。即使超参数设置不够好，TRPO 的表现也不会太差。相比之下，策略梯度算法对超参数更敏感。

TRPO 算法真正实现起来并不容易，主要难点在于第二步——**最大化**。不建议读者自己去实现 TRPO。

9.2 熵正则 (Entropy Regularization)

策略学习的目的是学出一个策略网络 $\pi(a|s; \theta)$ 用于控制智能体。每当智能体观测到当前状态 s ，策略网络输出一个概率分布，智能体依据概率分布抽样一个动作，并执行这个动作。举个例子，在超级玛丽游戏中，动作空间是 $\mathcal{A} = \{\text{左}, \text{右}, \text{上}\}$ 。基于当前状态 s ，策略网络的输出是

$$p_1 = \pi(\text{左} | s; \theta) = 0.3,$$

$$p_2 = \pi(\text{右} | s; \theta) = 0.4,$$

$$p_3 = \pi(\text{上} | s; \theta) = 0.3.$$

那么超级玛丽做的动作可能是左、右、上三者中的任何一个，概率分别是 0.3, 0.4, 0.3。在策略网络看来，每个动作都差不多，策略网络不知道哪个最好。假如对于任何状态 s ，策略网络的输出总是接近均匀分布，那么就意味着策略网络优柔寡断，没有判断力。我们不喜欢这种策略网络。相反，我们希望策略网络总是能够充满自信地做出决策，告诉智能体该做什么动作。举个例子，我们希望策略网络的输出是这样的概率分布

$$p_1 = \pi(\text{左} | s; \theta) = 0.05,$$

$$p_2 = \pi(\text{右} | s; \theta) = 0.83,$$

$$p_3 = \pi(\text{上} | s; \theta) = 0.12.$$

这个概率分布很明确地告诉智能体应该执行“向右”的动作，这个动作明显优于其余动作。可以用熵 (Entropy) 来衡量分布是否均匀。对于上述离散概率分布 $\mathbf{p} = [p_1, p_2, p_3]$ ，熵等于

$$\text{Entropy}(\mathbf{p}) = - \sum_{i=1}^3 p_i \cdot \ln p_i.$$

我们希望策略网络输出的概率分布有较小的熵；见图 9.3 的解释。

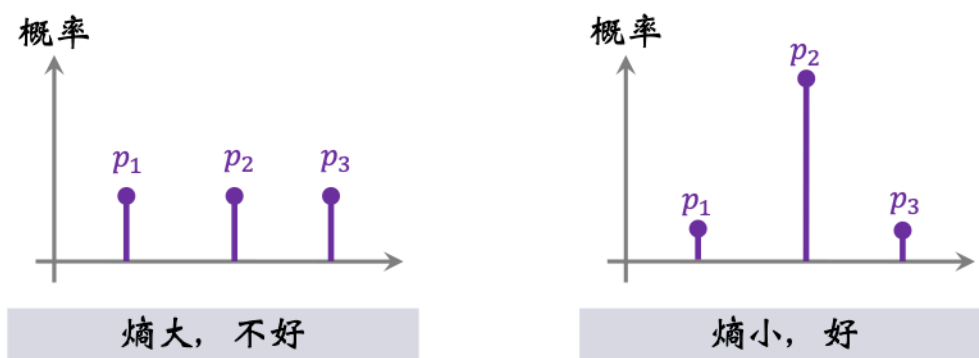


图 9.3: 两张图中分别描述两个离散概率分布。左边的概率密度很均匀，这种情况熵很大。右边的概率密度集中在 p_2 上，这种情况的熵较小。我们希望策略网络输出的概率值是右边的情况，即比较确信自己的判断。

策略学习中的熵正则：既然我们希望策略网络输出的概率分布有较小的熵，我们就不妨把熵作为正则项，放到策略学习的目标函数中。策略网络的输出是维度等于 $|\mathcal{A}|$ 的

向量，它表示定义在动作空间上的离散概率分布。这个概率分布的熵定义为：

$$H(s; \theta) \triangleq \text{Entropy} \left[\pi(\cdot | s; \theta) \right] = - \sum_{a \in \mathcal{A}} \pi(a | s; \theta) \cdot \ln \pi(a | s; \theta). \quad (9.11)$$

熵 $H(s; \theta)$ 只依赖于状态 s 与策略网络参数 θ 。我们希望对于大多数的状态 s ，熵都会比较小，也就是让 $\mathbb{E}_S[H(S; \theta)]$ 比较小。

回忆一下， $V_\pi(s)$ 是状态价值函数，衡量在状态 s 的情况下，策略网络 π 表现的好坏程度。策略学习的目标函数是 $J(\theta) = \mathbb{E}_S[V_\pi(S)]$ 。策略学习的目的是寻找参数 θ 使得 $J(\theta)$ 最大化。同时，我们还希望让熵最小化，所以把熵作为正则项，放到目标函数里。使用熵正则的策略学习可以写作这样的最大化问题：

$$\max_{\theta} J(\theta) - \lambda \cdot \mathbb{E}_S[H(S; \theta)]. \quad (9.12)$$

此处的 λ 是个超参数，需要手动调。

优化：带熵正则的最大化问题 (9.12) 可以用各种方法求解，比如策略梯度方法（包括 REINFORCE 和 Actor-Critic）、TRPO 等。此处只讲解策略梯度方法。公式 (9.12) 中目标函数关于 θ 的梯度是：

$$g(\theta) \triangleq \nabla_{\theta} [J(\theta) - \lambda \cdot \mathbb{E}_S[H(S; \theta)]].$$

观测到状态 s ，按照策略网络做随机抽样，得到动作 $a \sim \pi(\cdot | s; \theta)$ 。那么

$$\tilde{g}(s, a; \theta) \triangleq [Q_\pi(s, a) + \lambda \cdot \ln \pi(a | s; \theta) + \lambda] \cdot \nabla_{\theta} \ln \pi(a | s; \theta)$$

是梯度 $g(\theta)$ 的无偏估计（见定理 9.2）。因此可以用 $\tilde{g}(s, a; \theta)$ 更新策略网络的参数：

$$\theta \leftarrow \theta + \beta \cdot \tilde{g}(s, a; \theta).$$

此处的 β 是学习率。

定理 9.2. 带熵正则的策略梯度

$$\nabla_{\theta} [J(\theta) - \lambda \cdot \mathbb{E}_S[H(S; \theta)]] = \mathbb{E}_S [\mathbb{E}_{A \sim \pi(\cdot | s; \theta)} [\tilde{g}(S, A; \theta)]].$$

证明 首先推导熵 $H(S; \theta)$ 关于 θ 的梯度。由公式 (9.11) 中 $H(S; \theta)$ 的定义可得

$$\begin{aligned} \frac{\partial H(s; \theta)}{\partial \theta} &= - \sum_{a \in \mathcal{A}} \frac{\partial \pi(a | s; \theta) \cdot \ln \pi(a | s; \theta)}{\partial \theta} \\ &= - \sum_{a \in \mathcal{A}} \left[\ln \pi(a | s; \theta) \cdot \frac{\partial \pi(a | s; \theta)}{\partial \theta} + \pi(a | s; \theta) \cdot \frac{\partial \ln \pi(a | s; \theta)}{\partial \theta} \right]. \end{aligned}$$

第二个等式由链式法则得到。由于 $\frac{\partial \pi(a | s; \theta)}{\partial \theta} = \pi(a | s; \theta) \cdot \frac{\partial \ln \pi(a | s; \theta)}{\partial \theta}$ ，上面的公式可以写

成：

$$\begin{aligned}
\frac{\partial H(s; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} &= - \sum_{a \in \mathcal{A}} \left[\ln \pi(a|s; \boldsymbol{\theta}) \cdot \pi(a|s; \boldsymbol{\theta}) \cdot \frac{\partial \ln \pi(a|s; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} + \pi(a|s; \boldsymbol{\theta}) \cdot \frac{\partial \ln \pi(a|s; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \right] \\
&= - \sum_{a \in \mathcal{A}} \pi(a|s; \boldsymbol{\theta}) \cdot \left[\ln \pi(a|s; \boldsymbol{\theta}) + 1 \right] \cdot \frac{\partial \ln \pi(a|s; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \\
&= - \mathbb{E}_{A \sim \pi(\cdot|s; \boldsymbol{\theta})} \left[\left[\ln \pi(A|s; \boldsymbol{\theta}) + 1 \right] \cdot \frac{\partial \ln \pi(A|s; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \right]. \tag{9.13}
\end{aligned}$$

应用第 7 章推导的策略梯度定理，可以把 $J(\boldsymbol{\theta})$ 关于 $\boldsymbol{\theta}$ 的梯度写作

$$\frac{\partial J(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = \mathbb{E}_S \left\{ \mathbb{E}_{A \sim \pi(\cdot|S; \boldsymbol{\theta})} \left[Q_\pi(S, A) \cdot \frac{\partial \ln \pi(A|S; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \right] \right\}. \tag{9.14}$$

由公式 (9.13) 与 (9.14) 可得：

$$\begin{aligned}
&\frac{\partial}{\partial \boldsymbol{\theta}} \left[J(\boldsymbol{\theta}) - \lambda \cdot \mathbb{E}_S [H(S; \boldsymbol{\theta})] \right] \\
&= \mathbb{E}_S \left\{ \mathbb{E}_{A \sim \pi(\cdot|S; \boldsymbol{\theta})} \left[\left(Q_\pi(S, A) + \lambda \cdot \ln \pi(A|S; \boldsymbol{\theta}) + \lambda \right) \cdot \frac{\partial \ln \pi(A|S; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \right] \right\} \\
&= \mathbb{E}_S \left\{ \mathbb{E}_{A \sim \pi(\cdot|S; \boldsymbol{\theta})} [\tilde{g}(S, A; \boldsymbol{\theta})] \right\}.
\end{aligned}$$

上面第二个等式由 \tilde{g} 的定义得到。 □

相关文献

TRPO 由 John Schulman 等学者在 2015 年提出 [1]。TRPO 是置信域方法在强化学习中的成功应用。置信域是经典的数值优化算法，对此感兴趣的读者可以阅读这些教材：[2-3]。TRPO 每一轮循环都需要求解带约束的最大化问题；这类问题的求解可以参考这些教材：[4-5]。

参考文献

- [1] SCHULMAN J, LEVINE S, ABBEEL P, et al. Trust region policy optimization[C]//International Conference on Machine Learning (ICML). [S.l.: s.n.], 2015.
- [2] NOCEDAL J, WRIGHT S. Numerical optimization[M]. [S.l.]: Springer Science & Business Media, 2006.
- [3] CONN A R, GOULD N I, TOINT P L. Trust region methods[M]. [S.l.]: SIAM, 2000.
- [4] BERTSEKAS D P. Constrained optimization and lagrange multiplier methods[M]. [S.l.]: Academic press, 2014.
- [5] BOYD S, BOYD S P, VANDENBERGHE L. Convex optimization[M]. [S.l.]: Cambridge university press, 2004.
- [6] SILVER D, LEVER G, HEES N, et al. Deterministic policy gradient algorithms[C]//International Conference on Machine Learning (ICML). [S.l.: s.n.], 2014.
- [7] LILLICRAP T P, HUNT J J, PRITZEL A, et al. Continuous control with deep reinforcement learning.[C]//International Conference on Learning Representations (ICLR). [S.l.: s.n.], 2016.
- [8] HAFNER R, RIEDMILLER M. Reinforcement learning in feedback control[J]. Machine learning, 2011, 84 (1-2): 137-169.
- [9] PROKHOROV D V, WUNSCH D C. Adaptive critic designs[J]. IEEE transactions on Neural Networks, 1997, 8(5): 997-1007.