传统增强学习(RL)与变分方法的本质等价 性分析

作者: GaoZheng日期: 2025-03-18

• 版本: v1.0.0

增强学习(Reinforcement Learning, RL)与变分方法在本质上具有数学上的等价性,但它们在表达方式、计算方法和适用领域上有所不同。强化学习的优化过程可以视为一个基于期望泛函优化的变分问题,许多现代强化学习算法已直接使用变分方法进行推导和优化。

1. 变分方法的基本框架

变分方法 (Variational Methods) 用于求解泛函极值问题,即:

$$\delta S = 0$$

其中 $S[\pi]$ 是某种目标泛函,例如:

• 经典力学中的作用量:

$$S[q] = \int L(q,\dot{q},t) dt$$

• 量子力学中的基态能量:

$$E = \min_{\psi} rac{\langle \psi | H | \psi
angle}{\langle \psi | \psi
angle}$$

• 机器学习中的最优化问题:

$$\theta^* = \arg\min_{\theta} J(\theta)$$

在增强学习的背景下,变分方法可视为**策略优化过程的理论基础**,即寻找使得策略性能最优的泛函路 径。

2. 传统增强学习 (RL) 的优化过程

强化学习的核心目标是最大化累积回报:

$$J(\pi) = \mathbb{E}_{ au \sim \pi} \left[\sum_{t=0}^T \gamma^t R(s_t, a_t)
ight]$$

其中:

- $\pi(a|s)$ 是策略 (policy) , 表示在状态 s 下采取动作 a 的概率。
- R(s,a) 是奖励函数。
- γ 是折扣因子。
- 轨迹 $\tau = (s_0, a_0, s_1, a_1, \dots)$ 由策略 π 生成。

强化学习的目标是找到最优策略:

$$\pi^* = rg \max_{\pi} J(\pi)$$

这一优化目标本质上就是一个**变分优化问题**,其中**策略 π 是优化变量**,而目标函数是累积奖励。

3. RL 与变分方法的等价性

强化学习的优化过程可以直接映射到变分方法:

- 策略梯度方法 (Policy Gradient)
 - 通过梯度上升优化策略:

$$abla_{ heta} J(\pi_{ heta}) = \mathbb{E}\left[
abla_{ heta} \log \pi_{ heta}(a|s) Q^{\pi}(s,a)
ight]$$

- 。 这一公式可以看作是**基于概率测度的变分优化**,其中 $Q^{\pi}(s,a)$ 作为目标泛函,策略 π_{θ} 作为 变分参数。
- 基于变分推断的RL
 - 。 现代强化学习中, 许多方法已经直接采用**变分推断 (Variational Inference, VI) **方法:

$$\pi^* = rg \max_{\pi} \mathbb{E}_{ au \sim \pi} \left[R(au)
ight] - D_{ ext{KL}}(\pi || \pi_{ ext{prior}})$$

。 这表明 RL 在本质上已经转化为了一个变分优化问题。

• 贝叶斯RL: 从变分自由能到最优策略

。 在贝叶斯优化和信息论增强学习(Information-Theoretic RL)中,RL问题可以等价于最小化变分自由能(Variational Free Energy):

$$F = -\mathbb{E}_{\pi}[R(au)] + D_{\mathrm{KL}}(\pi||p(au))$$

。这一形式与统计物理中的变分方法完全一致。

4. RL 与变分方法的对比分析

对比维度	传统变分方法	传统增强学习(RL)
优化目标	目标泛函 $S[\pi]$	期望奖励函数 $J(\pi)$
优化变量	泛函路径 $q(t)$ 或概率密度 $p(x)$	策略 $\pi(a,s)$
求解方法	变分微分 $\delta S=0$	策略梯度,强化学习
信息论视角	变分自由能最小化	期望奖励最大化,KL 正则化
计算方式	解析求解 + 数值优化	采样 + 经验回放 + 变分近似

从该表格可以看出, RL 实际上是**一个特殊的变分优化过程**:

- 策略 π 对应于变分法中的泛函路径。
- 目标函数 $J(\pi)$ 对应于变分方法中的泛函极值。
- 变分方法通常使用梯度下降优化,而 RL 使用策略梯度和贝叶斯变分优化来优化策略。

5. RL 作为变分方法的泛化

尽管 RL 和变分方法在数学上等价,但 RL 在以下方面具有更广泛的适用性:

1. 强化学习可在高维环境下优化策略

• 变分方法通常适用于解析可解的问题,而 RL 可在高维、非线性、部分可观测环境下进行学习。

2. 强化学习可处理不确定性

- 变分方法通常假设问题是确定的,而 RL 通过探索-开发平衡在不确定环境下优化策略。
- 3. 强化学习可以利用经验回放进行训练

• 变分方法通常需要解析推导,而 RL 可通过数据采样和经验回放进行优化,使其更适用于大规模计算。

4. RL 可用于决策优化

• 变分方法通常用于能量最小化, 而 RL 可用于决策问题, 如机器人控制、金融交易等。

6. 结论

- 1. **传统增强学习和变分方法在数学上是等价的**,RL 是优化目标函数 $J(\pi)$ 的变分优化过程。
- 2. **许多现代RL方法已明确使用变分推导,如变分强化学习(Variational RL)**,其将RL等价于贝叶斯推断和变分自由能最小化问题。
- 3. RL 是变分方法的扩展,能够适应更复杂的环境,包括高维、部分可观测、不确定性优化等问题。

最终结论

- 传统变分方法 = 变分优化泛函极值问题
- 传统增强学习 = 变分方法在高维、不确定环境下的泛化
- GRL路径积分 = 变分方法 + 路径积分 + 泛范畴优化, 进一步拓展 RL 的计算能力

从这一角度来看,GRL 路径积分不仅是 RL 的扩展,更是变分方法与 RL 在泛范畴理论下的进一步统一,使其具备更强的计算能力和泛化性。

许可声明 (License)

Copyright (C) 2025 GaoZheng

本文档采用知识共享-署名-非商业性使用-禁止演绎 4.0 国际许可协议 (CC BY-NC-ND 4.0)进行许可。