

start.
PPO. 2025.10.13.

1): Pre-Request: Policy Gradient.

目标: 最大化在策略 θ 下的期望回报 ($R(\tau)$). return on τ

$$E(R(\tau))_{\tau \sim p_{\theta}} = \sum_{\tau} R(\tau) \cdot p_{\theta}(\tau).$$

故, 对其使用 梯度上升 \rightarrow maximize here.

$$\nabla_{\theta} E(R(\tau))_{\tau \sim p_{\theta}} = \nabla \sum_{\tau} R(\tau) \cdot p_{\theta}(\tau)$$

找 θ 变化的方向. $= \sum_{\tau} R(\tau) \nabla p_{\theta}(\tau)$. \swarrow $R(\tau)$ 与 θ 无关. \searrow $p_{\theta}(\tau)$ 难以直接得到.

Δ tip. 梯度 ∇ . $= \sum_{\tau} R(\tau) \cdot p_{\theta}(\tau) \cdot \frac{\nabla p_{\theta}(\tau)}{p_{\theta}(\tau)}$ (数学变换).

eg. $\nabla \log p_{\theta}(\text{action})$

指向 θ 空间内 (state) 处

使 $\log p_{\theta}(\text{action})$
增长最快的方向.

(需要回顾 ∇).

10.13.

由 Monte-Carlo \rightarrow 采样近似.

$$\begin{aligned} &= \frac{1}{N} \cdot \sum_{n=1}^N R(\tau^n) \cdot \frac{\nabla p_{\theta}(\tau^n)}{p_{\theta}(\tau^n)} \quad \text{又: } \nabla \log f(x) = \frac{\nabla f(x)}{f(x)} \\ &\quad \downarrow \text{类似求导} \\ &= \frac{1}{N} \cdot \sum_{n=1}^N R(\tau^n) \cdot \nabla \log p_{\theta}(\tau^n). \end{aligned}$$

eg: $T: \{s_1, a_1, s_2, a_2, \dots, s_n, a_n, s_{n+1}\}$.

接上.

$$\nabla E_{P(\mathcal{T})}(R(\mathcal{T})) \approx \frac{1}{N} \sum_{n=1}^N R(\mathcal{T}^n) \cdot \nabla \log P(\mathcal{T}^n) \quad (1).$$

需要处理.

这里认为, 状态转移仅由当前 s, a 决定.

即: $P(s_{t+1} | s_t, a_t)$. (即. 马尔可夫过程).

此时, $P(\mathcal{T})$ 可表示为: $P(\mathcal{T}) = p(s_1) \cdot \prod_{t=1}^T p(a_t | s_t) \cdot P(s_{t+1} | s_t, a_t)$.

$$P(\mathcal{T}) = \underbrace{p(s_1)}_{\text{初始 state 和 状态转移}} \cdot \underbrace{p(a_1 | s_1)}_{\text{由 Environment 决定}} \cdot \underbrace{p(s_2 | s_1, a_1)}_{\text{模型进行决策的过程}} \cdot \underbrace{p(a_2 | s_2)}_{\text{由 Environment 决定}} \cdots p(a_T | s_T) \cdot P(s_{T+1} | s_T, a_T)$$

初始 state 和 状态转移

由 Environment 决定

50 天关

$$\text{故有: } \nabla_{\theta} \log P(\mathcal{T}) = \underbrace{\nabla_{\theta} (\log p(s_1))}_{\theta \text{ 无关} \rightarrow 0} + \sum_{t=1}^{T-1} \underbrace{\nabla_{\theta} (\log p(a_t | s_t))}_{\theta \text{ 无关} \rightarrow 0} + \sum_{t=1}^T \underbrace{\nabla_{\theta} (\log P(s_{t+1} | s_t, a_t))}_{\theta \text{ 无关} \rightarrow 0}$$

$$= \sum_{t=1}^T \nabla_{\theta} \log p(a_t | s_t).$$

故. (1) 化为: $= \frac{1}{N} \sum_{n=1}^N R(\mathcal{T}^n) \cdot \sum_{t=1}^T \nabla_{\theta} \log p(a_t | s_t^n).$

理解:

当 $R(\mathcal{T}^n) > 0$,

算法会将 T^n 内所有 $p(a_t | s_t)$

的概率都加大.

(认为既然 $R(\mathcal{T}^n) > 0$, 则所有动作皆优).

显然, 要优化,

$$= \frac{1}{N} \sum_{n=1}^N \cdot \sum_{t=1}^T R(\mathcal{T}^n) \cdot \nabla_{\theta} \log p(a_t^n | s_t^n). \quad (2)$$

模型真正干的事: 由 state 做出决策.

相较于 $P(\mathcal{T}^n)$,

好求了许多.

现在的问题是:

直接用 $R(\mathcal{T}^n)$ 显然不合适.

如何优化?

10-13. end.

接上