

PPO 算法.

目标：训练一个 Policy π . 使其在所有 state 和 action 下，找一条运行 trajectory T . 其 reward $R(T)$ 最大。

$$\text{指导} : \max E(R(T))_{T \sim P_\theta(T)} = \sum_T R(T) P_\theta(T).$$

$$\begin{aligned} \Rightarrow \nabla E(R(T))_{T \sim P_\theta(T)} &= \sum_T R(T) \nabla P_\theta(T) \\ &= \sum_T R(T) \cdot P_\theta(T) \cdot \frac{\nabla P_\theta(T)}{P_\theta(T)} \\ &\approx \frac{1}{N} \sum_{n=1}^N R(T_n) \cdot \nabla \log P_\theta(T_n). \\ &= \frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{T_n} R(T_n) \nabla \log P_\theta(a_t^n | s_t^n). \end{aligned}$$

$$\Rightarrow \text{损失函数} : -\frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{T_n} R(T_n) \log P_\theta(a_t^n | s_t^n).$$

$R(T_n)$ 影响整个 action 链. 所有 action 被赋予相同 reward.

$$\Rightarrow R(T_n) \mapsto R_t^n = \sum_{t'=t}^{T_n} \gamma^{t'-t} \cdot r_{t'}^n$$

$$\Rightarrow -\frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{T_n} (R_t^n - B(s_t^n)) \cdot \log P_\theta(a_t^n | s_t^n) \Rightarrow B(s_t^n) \text{ 为平均价值. 给出}$$

训练中 $R_t^n - B(s_t^n)$ 用 $Q(s, a) - V_\theta(s)$ 代替.

多个 "好" action 中最 "好" 的

定义：① $Q(s, a)$ 为动作价值函数. \Rightarrow 判断 s 下 a 的价值.

② $V_\theta(s)$ 为状态价值函数.

③ $A_\theta(s, a)$ 为优势函数 $\Rightarrow A_\theta(s, a) = Q(s, a) - V_\theta(s)$.

关系： $Q_\theta(s_t, a) = r_t + \gamma \cdot V_\theta(s_{t+1})$. $V_\theta(s_{t+1}) \approx r_{t+1} + \gamma \cdot V_\theta(s_{t+2})$

目前损失函数： $-\frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{T_n} A_\theta(s_t^n, a_t^n) \cdot \log P_\theta(a_t^n | s_t^n)$

$A_\theta(s, a)$ 多步代入

$$\Rightarrow A_\theta^1(s_t, a) = Q_\theta(s_t, a) - V_\theta(s_t) = r_t + \gamma \cdot V_\theta(s_{t+1}) - V_\theta(s_t)$$

$$A_\theta^2(s_t, a) = r_t + \gamma \cdot r_{t+1} + \gamma^2 V_\theta(s_{t+2}) - V_\theta(s_t)$$

$$A_\theta^3(s_t, a) = r_t + \gamma \cdot r_{t+1} + \gamma^2 V_\theta(s_{t+2}) + \gamma^3 V_\theta(s_{t+3}) - V_\theta(s_t)$$

$$\text{定义: } \delta_t^V = r_t + \gamma V_\theta(s_{t+1}) - V_\theta(s_t)$$

$$\Rightarrow A_\theta^1(s_t, a) = \delta_t^V, \quad A_\theta^2(s_t, a) = \delta_t^V + \gamma \delta_{t+1}^V \dots \quad A_\theta^T(s_t, a) = \sum_{k=0}^{T-1} \gamma^k \delta_{t+k}^V$$

$$\text{定义 } A_\theta^{\text{GAE}}(s_t, a) = (1-\lambda)(A_\theta^1 + \gamma A_\theta^2 + \gamma^2 A_\theta^3 + \dots)$$

$$= \sum_{i=0}^{\infty} (\gamma \lambda)^i \delta_{t+i}^V$$

$$\Rightarrow \text{Loss} = -\frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{T_n} A_\theta^{\text{GAE}}(s_t^n, a_t^n) \cdot \log P_\theta(a_t^n | s_t^n).$$

其中, A_θ^{GAE} 中的 $V_\theta(s)$ 由神经网络拟合. label 为 $\sum_{t=t}^{T_n} \gamma^{t-t} r_t^n$.

重要性采样:

$$E(f(x))_{X \sim p(x)} = \sum_x p(x) f(x) = \sum_x q(x) \frac{p(x)}{q(x)} f(x) = E\left(\frac{p(x)}{q(x)} f(x)\right)_{X \sim q(x)}.$$

Off-Policy: 从 $q(x)$ 中采样, 训练 $p(x)$. 即可以从一个固定策略中采样.

用其从环境中取得的反馈训练新策略.

$$\text{Loss} = -\frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{T_n} A_{\theta'}^{\text{GAE}}(s_t^n, a_t^n) \cdot \frac{P_\theta(a_t^n | s_t^n)}{P_{\theta'}(a_t^n | s_t^n)} \cdot \log P_\theta(a_t^n | s_t^n).$$

$$\text{求导更新后等价于 } -\frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{T_n} A_{\theta'}^{\text{GAE}}(s_t^n, a_t^n) \frac{P_\theta(a_t^n | s_t^n)}{P_{\theta'}(a_t^n | s_t^n)} + \beta \text{KL}(P_\theta, P_{\theta'})$$

$$\text{截断策略: } -\frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{T_n} \min\left(A_{\theta'}^{\text{GAE}}(s_t^n, a_t^n) \frac{P_\theta(a_t^n | s_t^n)}{P_{\theta'}(a_t^n | s_t^n)}, \text{clip}\left(\frac{P_\theta(a_t^n | s_t^n)}{P_{\theta'}(a_t^n | s_t^n)}, 1-\epsilon, 1+\epsilon\right) A_{\theta'}^{\text{GAE}}(s_t^n, a_t^n)\right)$$

均限制 θ 与 θ' 分布差距不应过大.