

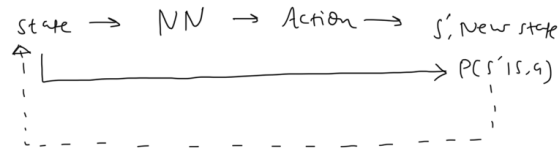
Introduction

Policy, $\pi_{\theta}(a|s)$, parametrized using ' θ '

→ We can use Deep NN to represent $\pi_{\theta}(a|s)$.

↓

' θ ' denotes the policy parameters which will be the weights of Deep Neural Network



Model-free: 1. Do not know $P(s_{t+1}|s_t, a_t)$
2. Do not know the initial state probability

Trajectory distribution

$$P_{\theta}(s_1, a_1, s_2, a_2, \dots, s_T, a_T) = P(s_1) \prod_{t=1}^T \pi_{\theta}(a_t|s_t) P(s_{t+1}|s_t, a_t)$$

$$P_{\theta}(\tau)$$

→ Trajectory distribution depends on the parameter ' θ '

REINFORCEMENT LEARNING OBJECTIVE

$$\theta^* = \underset{\theta}{\operatorname{argmax}} \mathbb{E}_{\tau \sim P_{\theta}(\tau)} \left[\sum_t r(s_t, a_t) \right] J(\theta)$$

Find the parameters θ which maximise the Expectation.

$$J(\theta) = \mathbb{E}_{\tau \sim P_{\theta}(\tau)} \left[\sum_{t=1}^T r(s_t, a_t) \right]$$

$$= \int P_{\theta}(\tau) r(\tau) d\tau$$

$$\nabla_{\theta} J(\theta) = \nabla_{\theta} \int P_{\theta}(\tau) r(\tau) d\tau = \int \nabla_{\theta} P_{\theta}(\tau) r(\tau) d\tau$$

differential operator is linear

→ Cannot differentiate w.r.t transition probabilities and initial state distribution is unknown.

→ What we do know are the trajectory samples and the rewards

$$\rightarrow P_{\theta}(\tau) \nabla_{\theta} \log P_{\theta}(\tau) = P_{\theta}(\tau) \frac{\nabla_{\theta} P_{\theta}(\tau)}{P_{\theta}(\tau)} = \nabla_{\theta} P_{\theta}(\tau)$$

$$\begin{aligned} \nabla_{\theta} J(\theta) &= \int P_{\theta}(\tau) \nabla_{\theta} \log P_{\theta}(\tau) r(\tau) d\tau \\ &= \mathbb{E}_{\tau \sim P_{\theta}(\tau)} \left[\nabla_{\theta} \log P_{\theta}(\tau) r(\tau) \right] \end{aligned}$$

sample from trajectory

$$\rightarrow P_{\theta}(s_1, a_1, \dots, s_T, a_T) = P(s_1) \prod_{t=1}^T \pi_{\theta}(a_t|s_t) P(s_{t+1}|s_t, a_t)$$

$$\rightarrow \log P_{\theta}(\tau) = \log P(s_1) + \sum_{t=1}^T \log \pi_{\theta}(a_t|s_t) + \log P(s_{t+1}|s_t, a_t)$$

$$\rightarrow \nabla_{\theta} \log P_{\theta}(\tau) = \underset{\text{independent of } \theta}{0} + \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(a_t|s_t) + \underset{0}{0}$$

$$\begin{aligned} \nabla_{\theta} J(\theta) &= \mathbb{E}_{\tau \sim P_{\theta}(\tau)} \left[\nabla_{\theta} \log P_{\theta}(\tau) r(\tau) \right] \\ &= \mathbb{E}_{\tau \sim P_{\theta}(\tau)} \left[\left(\sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(a_t|s_t) \right) \left(\sum_{t=1}^T r(s_t, a_t) \right) \right] \end{aligned}$$

$$\nabla_{\theta} J(\theta) = \frac{1}{N} \sum_{N=1}^N \left(\quad \right) \left(\quad \right)$$

$$\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\theta)$$