

π

$$y_t = r_t + \gamma \underbrace{V^\pi(s')}_{\Theta^\top \Phi(s')}$$

θ^*

$\min_{\theta} L = \mathbb{E} \left[\underbrace{V_{\theta}^{\pi}}_{\text{? States}} - y_t \right)^2$

Stationary state distribution after running policy π potentially for a long time

$$\lim_{t \rightarrow \infty} \theta_t = \theta^*$$

1997

$$J(\theta) = \mathbb{E}_{\tau \sim P_\theta(\tau)} [r(\tau)]$$

$P_\theta(\tau) \neq 0$ return ایسز، تاویل
احتمال یک ایسز

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

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

Under certain Continuity Cond.

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

$\nabla_\theta \log P_\theta(\tau)$

$$= \int \underbrace{\nabla_{\theta} \log P_{\theta}(\tau)}_{g(\tau)} r(\tau) P_{\theta}(\tau) d\tau$$

$$= \underset{\tau \sim P_{\theta}(\tau)}{E(g(\tau))} \approx \frac{1}{N} \sum_{i=1}^N \nabla_{\theta} \log \underbrace{P_{\theta}(\tau_i)}_{\text{blue underline}} r(\tau_i)$$

$$P_{\theta}(\tau_i) = P(s_0^{(i)}) P(s_1^{(i)} | s_0^{(i)}, a_0^{(i)}) \pi_{\theta}(a_0^{(i)} | s_0^{(i)}) \dots$$

$$\begin{aligned} \nabla_{\theta} \log P_{\theta}(\tau) &= \nabla_{\theta} \log P(s_0) + \nabla_{\theta} \log P(s_1 | s_0, a_0) \\ &\quad + \nabla_{\theta} \log \pi_{\theta}(a_0 | s_0) + \dots \end{aligned}$$

(green underline and text below the last term)

$$\nabla_{\theta} J(\theta) \equiv \frac{1}{N} \sum_{i=1}^N \left(\sum_{j=1}^T \nabla \log \pi(a_j^{(i)} | s_j^{(i)}) \right).$$

شیوه MLE وزن دار
هری return اون episode بیشتر
باشه در شبکه بیشتر اثر طره

$$\left(\sum_{j=1}^T r_j^{(i)} \right)$$

MC

τ_1

وارایش بالا

τ_2

①

②

بازی عقلی، ارزش نمونه بگیر
نمونه آپدیت به این صورت که احتمال درگیر نمونه خوب رو ببر بالا
مفهوم مسدود شدن یعنی احتمال return های خوب رو ببر بالا در مکن احتمال state های

در برابر پایین

τ_N

$$\theta^{(t+1)} \leftarrow \theta^{(t)} + \alpha \nabla_{\theta} J(\theta)$$