

π

$$y_t = r_t + \gamma \underbrace{V^\pi(s')}_{\underbrace{\Theta^T \Phi(s')}_T}$$

$$\theta^* = \min_{\theta} L = \mathbb{E} \left[\underbrace{V^{\pi}_{\theta}}_{\text{? States} \rightarrow \text{stationary state distribution after running policy } \pi} - y_t \right)^2$$

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

1997

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

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

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

Under certain Continuity Cond.

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

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

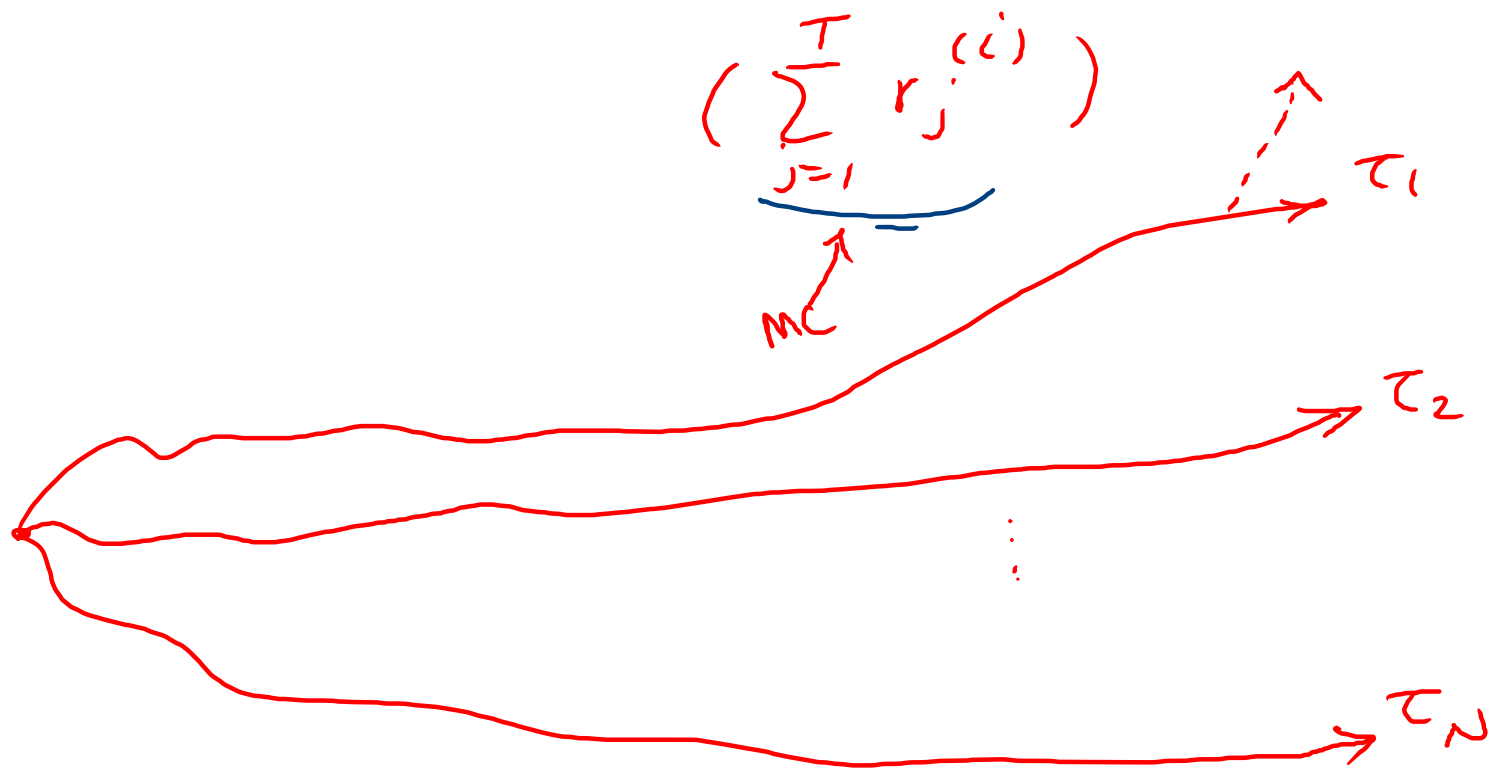
$$= \int \underbrace{\nabla_{\theta} \log P_{\theta}(\tau)}_{g(\tau)} r(\tau) P_{\theta}(\tau) d\tau$$

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

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

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

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



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