

$$Q(s,a;\theta)$$

$$\theta$$

$$Q^*(s,a)\approx$$

$$Q(s,a:$$

$$\theta)$$

$$\theta$$

$$L_i(\theta_i)=E\Big(r+\gamma\max_{a'}Q(s',a';\theta_{i-1})-Q(s,a:\theta_i)\Big)^2$$

$$\frac{s'}{\theta}\in$$

$$R_t^d$$

$$\pi_\theta(a|s)$$

$$\theta$$

$$J(\pi_\theta)$$

$$J(\pi_\theta)$$

$$J$$

$$J(\pi_\theta)=$$

$$E[R_t]$$

$$J$$

$$\theta_{t+1}=\theta_t+\alpha\nabla J(\theta_t)$$

$$\nabla \hat{J}(\theta_t)$$

$$J_{\theta_t}$$

$$\nabla J(\pi_\theta) \propto \sum_s \mu(s) \sum_a q_\pi(s,a) \nabla_\theta \pi_\theta(a|s)$$

$$\frac{\mu}{S}$$

$$\frac{s}{\pi_\theta}$$

$$?$$

$$\nabla J(\pi_\theta) = E_\pi \left[R_t \frac{\nabla_\theta \pi_\theta(a|S_t)}{\pi_\theta(a|S_t)} \right]$$

$$\left[R_t\frac{\nabla_\theta\pi_\theta(a|S_t)}{\pi_\theta(a|S_t)}\right]$$

$$\nabla J(\pi_\theta)$$

$$\theta$$

$$\theta_{t+1}=\theta_t+\alpha R_t\frac{\nabla_\theta\pi_\theta(a|S_t)}{\pi_\theta(a|S_t)}=\theta_t+\alpha R_t\nabla_\theta log\pi_\theta(a|S_t)$$

$$J(\pi\theta)$$

$$?$$

$$\theta$$

$$\nabla_\theta log\pi(a_t|s_t:$$

$$\theta)R_t$$

$$\nabla_\theta E[R_t]$$

$$b_t(s_t)$$

$$\bar{R}_t$$

$$\nabla_\theta log\pi(a_t|s_t;\theta)(R_t-$$

$$b_t(s_t))$$

$$b_t(s_t)\approx$$

$$V^\pi(s_t)$$

$$\pi$$

$$\bar{b}_t$$

$$b_t(s)=$$

$$\hat{V}_{\pi_{\theta_t}}(s)$$

$$R_t-$$

$$b_t(s)$$

$$a_t$$

$$s_t$$

$$\hat{A}(a_t,s_t)=$$

$$Q(a_t,s_t)-$$

$$V(s_t)$$

$$R_t$$

$$Q^\pi(a_t,s_t)$$

$$b_t$$

$$\hat{V}^\pi(s_t)$$

$$AdvantageActor-$$

$$Critic$$

$$A2C$$

$$AdvantageActor-$$

$$A2C$$