

$$\pi(a|s, \bar{\theta})$$

$$J(\bar{\theta})$$

$$\bar{\theta}_{t+1} = \bar{\theta}_t + \alpha \nabla_{\bar{\theta}} J(\bar{\theta}) \big|_{\bar{\theta} = \bar{\theta}_t}$$

$$\bar{\theta}^T x(s, a)$$

$$h(s, a; \bar{\theta}) \text{ — "score"}$$

$$\pi(a|s, \bar{\theta}) = \text{softmax}(h) = \frac{e^{h(s, a; \bar{\theta})}}{\sum_{a'} e^{h(s, a'; \bar{\theta})}}$$

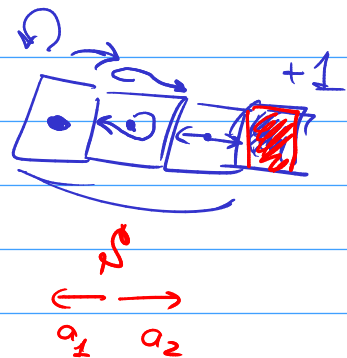
greedy

модель

1) Сход. к грейду сфас.

2)  $\pi$  и  $Q$  пропе ант, чем  $Q$

3) Модель сход. к стохаст. сфасеру



$$J(\theta) = V_{\pi_{\theta}}(s_0)$$

Policy gradient theorem

$$\nabla_{\theta} V_{\pi}(s) = \nabla_{\theta} \left[ \sum_a \pi_{\theta}(a|s) Q_{\pi_{\theta}}(s, a) \right] =$$

$$\delta = 1$$

$$= \sum_a \left[ \nabla_{\theta} \pi_{\theta}(a|s) \cdot Q_{\pi_{\theta}}(s, a) + \pi_{\theta}(a|s) \cdot \nabla_{\theta} Q_{\pi_{\theta}}(s, a) \right] =$$

$$= \sum_a \left[ \nabla_{\theta} \pi_{\theta}(a|s) \cdot Q_{\pi_{\theta}}(s, a) + \pi_{\theta}(a|s) \cdot \nabla_{\theta} \sum_{s'} p(s'|s, a) (r + V_{\pi}(s')) \right]$$

$$= \sum_a \left[ - \dots + \pi_{\theta}(a|s) \cdot \nabla_{\theta} \sum_{s'} p(s'|s, a) \cdot V_{\pi}(s') \right]$$

$$= \sum_a \left[ \nabla_{\theta} \pi_{\theta}(a|s) \cdot Q_{\pi_{\theta}}(s, a) + \pi_{\theta}(a|s) \cdot \sum_{s'} p(s'|s, a) \cdot \nabla_{\theta} V_{\pi}(s') \right]$$

$$\sum_a \pi_{\theta}(a|s) \sum_{s'} p(s'|s, a)$$

$$\sum_{a'} \left[ \nabla \pi(a'|s) \right]$$

$$= \sum_{s'} \left( \sum_{k=0}^{\infty} \Pr_{\pi} [s \rightarrow s' \text{ in } k \text{ steps}] \right) \sum_a \nabla_{\theta} \pi_{\theta}(a|s') \cdot Q_{\pi}(s', a)$$

$\propto \mu(s)$

$$\nabla_{\theta} J(\theta) = \nabla_{\theta} V_{\pi}(s_0) = \sum_s \left( \sum_{k=0}^{\infty} \Pr[s_0 \rightarrow s \text{ in } k \text{ steps}] \right) \sum_a \nabla_{\theta} \pi(a|s) Q(s, a)$$

symm. spec., prob. of  $s$   
 $\mu_{\pi}(s) = \text{avg. prob. spec. to state } s$

$$\nabla_{\theta} J(\theta) \propto \sum_s \mu_{\pi}(s) \cdot \sum_a Q_{\pi}(s, a) \cdot \nabla_{\theta} \pi_{\theta}(a|s) =$$

Policy  
grad  
thm

$$= \mathbb{E}_{\pi} \left[ \sum_a Q_{\pi}(S_t, a) \nabla_{\theta} \pi_{\theta}(a|S_t) \right]$$

1) „All-actions“ algorithm  $\pi: S_0, A_0, S_1, \underline{A_1} \rightarrow$

$$\bar{\theta}_{t+1} = \bar{\theta}_t + \alpha \sum_a \hat{Q}(S_t, a, \bar{w}) \cdot \nabla_{\theta} \pi_{\theta}(a|S_t)$$

2) REINFORCE (1992)

$$\nabla J(\theta) \propto \mathbb{E}_{\pi} \left[ \sum_a Q_{\pi}(S_t, a) \nabla_{\theta} \pi_{\theta}(a|S_t) \cdot \frac{\pi_{\theta}(a|S_t)}{\sum_a \pi_{\theta}(a|S_t)} \right] =$$

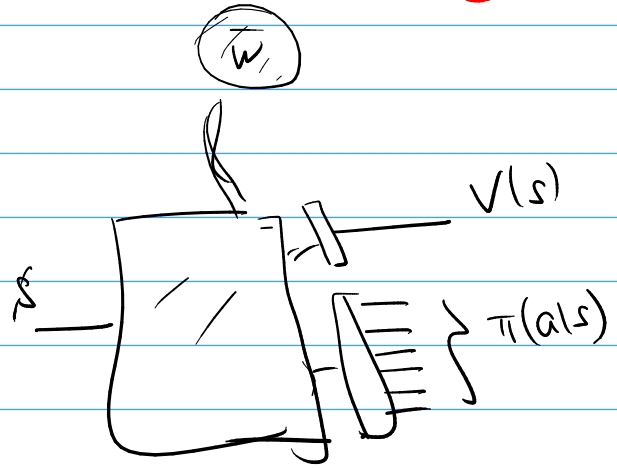
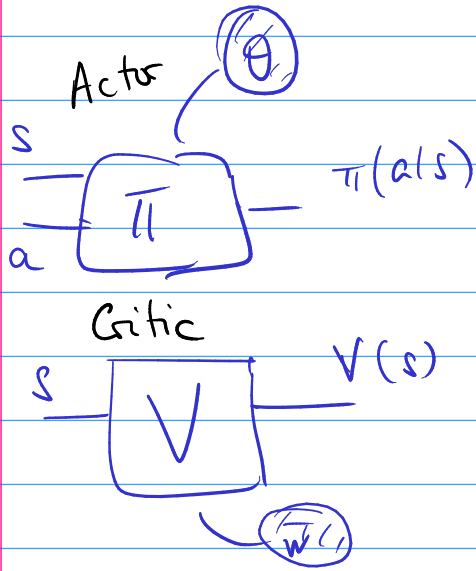
$$= \mathbb{E}_{\pi} \left[ \underbrace{Q_{\pi}(S_t, A_t)}_{\text{„}\mathbb{E}_{\pi}[G_t]\text{“}} \cdot \underbrace{\left( \frac{\nabla_{\theta} \pi_{\theta}(A_t|S_t)}{\pi_{\theta}(A_t|S_t)} \right)}_{\text{„}\nabla_{\theta} \ln \pi_{\theta}(A_t|S_t)\text{“}} \right] =$$

$$\boxed{\nabla J(\theta) \propto \mathbb{E}_{\pi} [G_t \cdot \nabla_{\theta} [\ln \pi_{\theta}(A_t|S_t)]]}$$

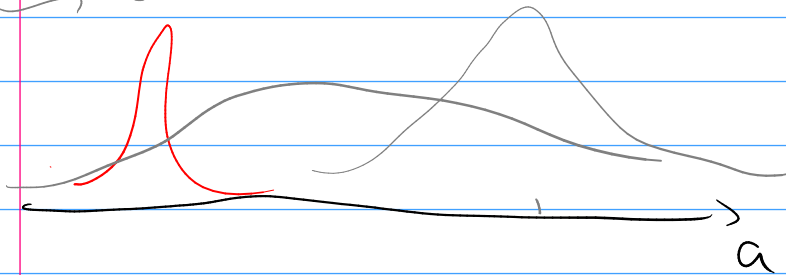
REINFORCE:  $\theta_{t+1} = \theta_t + \alpha \cdot G_t \cdot \nabla_{\theta} \ln \pi_{\theta}(A_t | S_t) =$

$\pi \sim \dots S_t, A_t, G_t, \dots$

$$= \theta_t + \alpha \cdot \frac{G_t}{\pi(A_t | S_t)} \cdot \nabla_{\theta} \pi(A_t | S_t)$$



$s' \rightsquigarrow a \sim \mathcal{N}(a | \mu_s, \sigma_s)$



$$\pi(a | s, \theta) = \frac{1}{\sqrt{2\pi} \cdot \sigma} e^{-\frac{1}{2\sigma^2} (a - \mu)^2}$$

$\mu(s, \bar{\theta})$

$\sigma(s, \bar{\theta})$

$$F(\bar{\theta}) \rightarrow \max$$

$$g(\bar{\theta}) \leq \delta$$

$$F(\bar{\theta}) - \frac{1}{\epsilon} g(\bar{\theta}) \rightarrow \max$$