

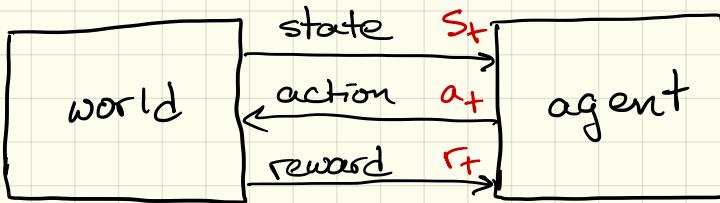
598 T2L

Fall 2018

T. Breitl

Policy Gradient (in class)

## MODEL



$$P(s_{t+1} | s_t, a_t)$$

"policy"

$$P(a_t | s_t ; \theta)$$

POLICY

$$\leftarrow \pi_\theta(a_t | s_t)$$

$$r_t = r(s_t, a_t)$$

sometimes called  $r_{t+1}$

GOAL

$$\tau = (s_1, a_1, \dots, s_T, a_T, s_{T+1}) \leftarrow \text{trajectory}$$

$$p(\tau; \theta) = p(s_1) p(a_1 | s_1; \theta) p(s_2 | s_1, a_1) \dots$$

$$= p(s_1) \prod_{t=1}^T p(a_t | s_t; \theta) p(s_{t+1} | s_t, a_t)$$

PAYOFF

$$\left\{ J(\theta) = E_{\tau \sim p(\tau; \theta)} \left[ \sum_{t=1}^T r(s_t, a_t) \right] \right.$$

$\uparrow$

MAXIMIZE

r( $\tau$ ) total reward

## POLICY GRADIENT

we want to find

$$\Theta^* = \arg \max_{\Theta} J(\Theta)$$

we can do that by gradient ascent

$$\Theta_{k+1} = \Theta_k + \alpha \boxed{\nabla_{\Theta} J(\Theta_k)}$$

$$J(\theta) = E_{\tau \sim P(\tau; \theta)} [r(\tau)]$$

$$\approx -\frac{1}{N} \sum_{i=1}^N r(\tau^i)$$

$\leftarrow \tau^1, \dots, \tau^N$  are generated by  $\theta$

$$= -\frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T r(s_t^i, a_t^i)$$

$$\frac{\partial J(\theta)}{\partial \theta} = -\frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \left( \frac{\partial r}{\partial s}(s_t^i, a_t^i) \frac{\partial s_t^i}{\partial \theta} + \frac{\partial r}{\partial a}(s_t^i, a_t^i) \frac{\partial a_t^i}{\partial \theta} \right)$$

$$r(s, a) = -(s - \mu)^2$$

$$\begin{aligned}
 J(\theta) &= E_{\tau \sim P(\tau; \theta)} [r(\tau)] \\
 &= \int_{\tau} r(\tau) p(\tau; \theta) d\tau
 \end{aligned}$$

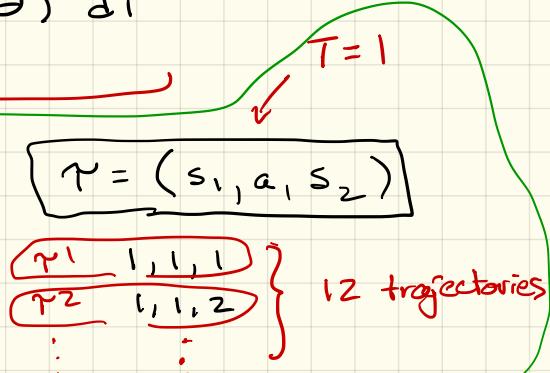
probability of sampling  $\tau$  under  $\theta$   
 $r(\tau)$  total reward for  $\tau$

integrate over the space of all trajectories

$$\nabla_{\theta} J(\theta) = \int_{\tau} r(\tau) \nabla_{\theta} p(\tau; \theta) d\tau$$

$\lambda = \{1, 2\}$   
 $S = \{1, 2, 3\}$

$$\frac{1}{12} \sum_{i=1}^{12} r(\tau^i) \nabla_{\theta} p(\tau^i; \theta)$$



## QUEUE

- infinite horizon ?
- discounted reward ?

$$\nabla_{\theta} J(\theta) = \int_{\gamma} r(\gamma) \ n_{\theta} p(\gamma; \theta) d\gamma$$

This does not have the form

$$\int_{\gamma} f(\gamma) p(\gamma; \theta) d\gamma = E_{\gamma \sim p(\gamma; \theta)} [f(\gamma)]$$

If it did, then we could approximate

$$E_{\gamma \sim p(\gamma; \theta)} [f(\gamma)] = \frac{1}{N} \sum_{i=1}^N f(\gamma_i)$$

$\gamma^1, \dots, \gamma^N$  are generated by  $\theta$

$$\nabla_{\Theta} J(\Theta) = \int_{\mathcal{R}} r(\tau) \underbrace{\delta_{\Theta} p(\tau; \Theta)}_{\text{red}} d\tau$$

$$p(r; \Theta) = p(s_1, a_1, \dots, s_T, a_T, s_{T+1})$$

$$= p(s_1) \prod_{t=1}^T p(a_t | s_t; \Theta) p(s_{t+1} | s_t, a_t)$$

$$\nabla_{\Theta} p(r; \Theta) = p(s_1) \sum_{t=1}^T \left( \prod_{k=1}^{t-1} p(a_k | s_k; \Theta) p(s_{k+1} | s_k, a_k) \right.$$

$$\left. \cdot \left( \nabla_{\Theta} p(a_t | s_t; \Theta) \right) p(s_{t+1} | s_t, a_t) \right)$$

$$\log p(\tau; \theta) = \log \left( p(s_1) \prod_{t=1}^T p(a_t | s_t; \theta) p(s_{t+1} | s_t, a_t) \right)$$

$$= \log p(s_1) + \sum_{t=1}^T (\log p(a_t | s_t; \theta) + \log p(s_{t+1} | s_t, a_t))$$

$$\nabla_{\theta} \log p(\tau; \theta) = \sum_{t=1}^T \nabla_{\theta} \log p(a_t | s_t; \theta)$$

\_\_\_\_\_

$$= \frac{1}{p(\tau; \theta)} \nabla_{\theta} p(\tau; \theta)$$

$$\nabla_{\theta} p(\tau; \theta) = p(\tau; \theta) \nabla_{\theta} \log p(\tau; \theta)$$

$$\nabla_{\theta} J(\theta) = \int_{\mathcal{T}} r(\tau) \nabla_{\theta} p(\tau; \theta) d\tau$$

$$= \int_{\mathcal{T}} (r(\tau) \nabla_{\theta} \log p(\tau; \theta)) p(\tau; \theta) d\tau$$

$$J(\theta) = E_{\gamma \sim p(\gamma; \theta)} [r(\gamma)]$$

$$= \int_{\gamma} r(\gamma) p(\gamma; \theta) d\gamma$$

$$\nabla_{\theta} J(\theta) = \int_{\gamma} r(\gamma) \boxed{\nabla_{\theta} p(\gamma; \theta)} d\gamma$$

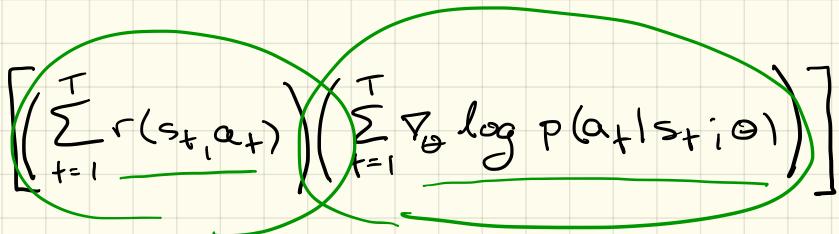
$$= \int_{\gamma} r(\gamma) \boxed{(\nabla_{\theta} \log p(\gamma; \theta)) p(\gamma; \theta)} d\gamma$$

$$\sum_{t=1}^T \nabla_{\theta} \log p(a_t | s_t; \theta)$$

$$= E_{\gamma \sim p(\gamma; \theta)} [r(\gamma) \nabla_{\theta} \log p(\gamma; \theta)]$$

## PAYOFF GRADIENT

$$\nabla_{\theta} J(\theta) = E_{\tau \sim p(s_t; \theta)} \left[ \left( \sum_{t=1}^T r(s_t, a_t) \right) \left( \sum_{t=1}^T \nabla_{\theta} \log p(a_t | s_t; \theta) \right) \right]$$



REINFORCE algorithm

① initialize  $\theta_0$

② sample  $\tau^1, \dots, \tau^N$  by running the policy  $p(a_t | s_t; \theta_0)$

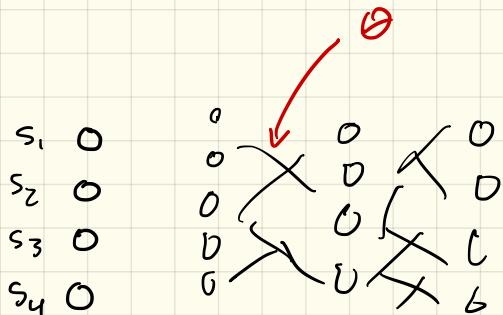
③ estimate  $\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_k \left( \left( \sum_{t=1}^T r(s_t^k, a_t^k) \right) \left( \sum_{t=1}^T \nabla_{\theta} \log p(a_t^k | s_t^k; \theta_k) \right) \right)$

④ update  $\theta_{k+1} = \theta_k + \alpha \nabla_{\theta} J(\theta_k)$

repeat  
until  
convergence

$$p(a|s_t; \theta)$$

(e.g., for the acrobot)



$$\begin{array}{ll} 0 & y_1 \quad (-M) \\ 0 & y_2 \quad (0) \\ 0 & y_3 \quad (M) \end{array}$$

$$p(a|s; \theta) = \begin{cases} y_1 & \text{if } a=a_1 \\ y_2 & \text{if } a=a_2 \\ y_3 & \text{if } a=a_3 \end{cases}$$

$$p(a|s; \theta)$$

$$p(a|s; \theta) = \begin{cases} \frac{e^{y_1}}{\sum_{i=1}^3 e^{y_i}} & \text{if } a=a_1 \\ \frac{e^{y_2}}{\cdot} & \text{if } a=a_2 \\ \frac{e^{y_3}}{\cdot} & \text{if } a=a_3 \end{cases}$$

$$\begin{cases} 1 & a=a_i \\ 0 & a \neq a_i \end{cases}$$

$$\log \left( \frac{e^{y_1}}{\sum_{i=1}^3 e^{y_i}} \right)$$
$$= y_1 - \log \left( \sum_{i=1}^3 e^{y_i} \right)$$



cross-entropy (basically)

Example. Gaussian policy

$$\begin{array}{ccccccccc} s_1 & 0 & & 0 & & 0 \\ & X & & & & & \\ s_2 & 0 & & 0 & X & 0 & X \\ & X & & & X & 0 & \nearrow \\ s_3 & 0 & & 0 & X & 0 & \nearrow \\ & X & & 0 & & & \\ s_4 & 0 & & X & 0 & X & \end{array}$$

0 u  
↓

$$p(a|s; \theta) = \mathcal{N} e^{-\frac{1}{2}(a - f(s; \theta))^T \Sigma^{-1} (a - f(s; \theta))}$$

$$\log p(a|s; \theta) = \log \mathcal{N} - \frac{1}{2}(a - f(s; \theta))^T \Sigma^{-1} (a - f(s; \theta))$$

$$\nabla_{\theta} \log p(a|s; \theta) = - \underbrace{\frac{\partial f^T}{\partial \theta}}_{\Sigma^{-1}} (a - f(s; \theta))$$

OFF-POLICY

importance sampling

$$\begin{aligned} E_{x \sim p(x)} [f(x)] &= \int_x p(x) f(x) dx \\ &= \int_x q_b(x) \left( \frac{p(x)}{q_b(x)} \right) f(x) dx \\ &= E_{x \sim q_b(x)} \left[ \left( \frac{p(x)}{q_b(x)} \right) f(x) \right] \end{aligned}$$

$$E_{\tau \sim p(\tau; \theta')} \left[ \nabla_{\theta} \log p(a_t | s_t; \theta') r(\tau) \right]$$

$$= E_{\tau \sim p(\tau; \theta)} \left[ \frac{p(\tau; \theta')}{p(\tau; \theta)} \nabla_{\theta} \log p(a_t | s_t; \theta') r(\tau) \right]$$

this allows us to reuse old data  
and, maybe, to get better stability

$$\frac{p(r; \Theta')}{p(T; \Theta)} = \frac{\cancel{p(s_1)} \prod_{t=1}^T p(a_t | s_t; \Theta') p(s_{t+1} | s_t, a_t)}{\cancel{p(s_1)} \prod_{t=1}^T p(a_t | s_t; \Theta) \cancel{p(s_{t+1} | s_t, a_t)}}$$

$$= \prod_{t=1}^T \frac{p(a_t | s_t; \Theta')}{p(a_t | s_t; \Theta)}$$