

598 RL

Fall 2018

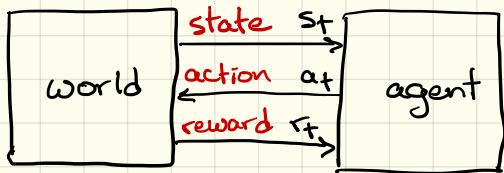
T. Breit

and M. West,  
personal correspondence

## Policy Gradient

(see also Levine, CS294, Lecture 5, Fall 2018)

## MODEL



$p(s_{t+1} | s_t, a_t)$  = probability that the next state is  $s_{t+1}$   
given that the state is  $s_t$  and the agent  
takes action  $a_t$

often called  $\pi_\theta(a_t | s_t)$

$p(a_t | s_t; \theta)$  = probability that the agent takes action  $a_t$   
given that the state is  $s_t$

POLICY ↗

↑  
it is parameterized by  $\theta$

$r_t = r(s_t, a_t)$  = reward for taking action  $a_t$  at state  $s_t$

## GOAL

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

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

↑ probability of generating this trajectory with a given policy

PAYOFF

WE WANT TO  
MAXIMIZE THIS

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

$r(\tau)$  total reward

take expectation wrt the distribution  $p(\tau; \theta)$   
on the space of trajectories

## METHOD - POLICY GRADIENT

we want to find

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

we can do so by gradient ascent

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

↑  
learning rate (step size)

← NEED TO COMPUTE THIS

## FIRST TRY

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

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

$\tau^1, \dots, \tau^N$  generated with  $\theta$

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

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_{\theta} r(s_t^i, a_t^i)$$

$$\underbrace{\frac{\partial r}{\partial s}(s_t^i, a_t^i) \frac{\partial s_t^i}{\partial \theta}} + \frac{\partial r}{\partial a}(s_t^i, a_t^i) \frac{\partial a_t^i}{\partial \theta}$$

↑  
this makes no sense at all  
(and can't be computed)

## SECOND TRY

$$J(\theta) = E_{\tau \sim P(\tau; \theta)} [r(\tau)]$$
$$= \int_{\text{space of trajectories}} r(\tau) P(\tau; \theta) d\tau$$

↑ total reward for this trajectory

probability of generating this trajectory

$$\nabla_{\theta} J(\theta) = \int_{\tau} \frac{r(\tau)}{\uparrow} \frac{\nabla_{\theta} P(\tau; \theta)}{\uparrow} d\tau$$

how good is  $\tau$ ?

what change in  $\theta$  would make the probability of  $\tau$  increase the fastest?

PROBLEM #1

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

This does not have the form

$$\int_{\tau} f(\tau) p(\tau; \theta) d\tau$$

so it can't be approximated by sampling  $\tau$  with  $\theta$  as

$$\frac{1}{N} \sum_{i=1}^N f(\tau^i),$$

Instead,  $\tau$  must be sampled uniformly, and so  $r(\tau)$  is likely to be very small.

PROBLEM #2

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

$$\begin{aligned} p(\tau; \theta) &= p(s_1, a_1, \dots, s_T, a_T, s_{T+1}; \theta) \\ &= p(s_1) \prod_{t=1}^T p(a_t | s_t; \theta) p(s_{t+1} | s_t, a_t) \end{aligned}$$

$$\begin{aligned} \nabla_{\theta} p(\tau; \theta) &= p(s_1) \sum_{t=1}^T \left( \left( \prod_{k=1}^T p(a_k | s_k; \theta) p(s_{k+1} | s_k, a_k) \right) \right. \\ &\quad \left. \cdot p(a_t | s_t; \theta) p(s_{t+1} | s_t, a_t) \right) \end{aligned}$$

this is a long product of probabilities, which is also likely to be very small (and numerically bad)

if only we were working with log-probabilities ...

$$\begin{aligned}\log p(r; \theta) &= \log \left( p(s_1) \prod_{t=1}^T p(a_t | s_t; \theta) p(s_{t+1} | s_t, a_t) \right) \\ &= \log p(s_1) + \sum_{t=1}^T \left( \log p(a_t | s_t; \theta) + \log p(s_{t+1} | s_t, a_t) \right)\end{aligned}$$

then ...

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



this is a sum and not a product  
and so will be much better behaved

↑  
it also doesn't need  $p(s_{t+1} | s_t, a_t)$  !!

BEHOLD!

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

what we want to find

what we are supposed to find



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



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

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

$$= \mathbb{E}_{\tau \sim p(\tau; \theta)} \left[ r(\tau) \nabla_{\theta} \log p(\tau; \theta) \right]$$

↑  
can be computed by sampling  $\tau$  with  $\theta$

## THE PAYOFF GRADIENT

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

$$\approx \frac{1}{N} \sum_{i=1}^N \left( \sum_{t=1}^T \nabla_{\theta} \log p(a_t^i | s_t^i; \theta) \right) \left( \sum_{t=1}^T r(s_t^i, a_t^i) \right)$$

$\uparrow$   
 $r^i = (s_1^i, a_1^i, \dots, s_T^i, a_T^i, s_{T+1}^i)$

is generated with  $\theta$

## REINFORCE algorithm

- ① sample  $\tau^1, \dots, \tau^N$  by running the policy  $p(a_t | s_t; \theta)$
  - ② estimate  $\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \left( \sum_{t=1}^T \nabla_{\theta} \log p(a_t^i | s_t^i; \theta) \right) \left( \sum_{t=1}^T r(s_t^i, a_t^i) \right)$
  - ③ update  $\theta_{k+1} = \theta_k + \alpha \nabla_{\theta} J(\theta_k)$
- repeat until convergence

↑  
how to compute this?

Example 1. Finite state and action space (tabular policy)

$\theta_{as}$  = weight of action  $a$  in state  $s$

$$p(a|s; \theta) = \frac{e^{\theta_{as}}}{\sum_{a'=1}^{na} e^{\theta_{a's}}}$$

softmax - exponentiate to make sure it is positive, then normalize

$$\log p(a|s; \theta) = \theta_{as} - \log \left( \sum_{a'=1}^{na} e^{\theta_{a's}} \right)$$

$$\begin{aligned} \frac{\partial}{\partial \bar{\theta}_{\bar{a}\bar{s}}} \log p(a|s; \theta) &= \delta_{\bar{a}a} \delta_{\bar{s}s} - \left( \sum_{a'=1}^{na} e^{\theta_{a's}} \right)^{-1} e^{\theta_{\bar{a}s}} \delta_{\bar{s}s} \\ &= \delta_{\bar{s}s} (\delta_{\bar{a}a} - p(\bar{a}|s; \theta)) \end{aligned}$$

Example 2. Continuous action space (Gaussian policy)

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

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

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

↓

$$\nabla_{\theta} \log p(a|s; \theta) = - \underbrace{\frac{\partial f^T}{\partial \theta} \Sigma^{-1}}_{\text{remember, it is easy to compute this by backpropagation if } f \text{ is a neural net}} (a - f(s; \theta))$$

remember, it is easy to compute  
this by backpropagation if  $f$  is  
a neural net

## VARIANCE (Lall, E207)

if  $x$  is a real-valued random variable, then

$$P(|x - E[x]| \geq a) \leq \frac{1}{a^2} \text{var}[x] \quad \} \text{Chebyshev inequality}$$

if  $x_1, \dots, x_N$  are real-valued and IID with mean  $\mu$  and variance  $\sigma^2$ , then

$$P\left(\left|\frac{1}{N} \sum_{i=1}^N x_i - \mu\right| \leq \epsilon\right) \geq 1 - \frac{\sigma^2}{N\epsilon^2} \quad \} \text{law of large numbers}$$

the confidence width that can be achieved at probability  $P_{\text{conf}}$  is

$$\epsilon = \sqrt{\frac{\sigma^2}{N(1-P_{\text{conf}})}}$$

$$\leftarrow N = \frac{\sigma^2}{\epsilon^2(1-P_{\text{conf}})}$$

## CAUSALITY

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

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

$\downarrow$  see next page

$$= \sum_{t=1}^T E_{\underbrace{(s_1, a_1, \dots, s_t, a_t, s_{t+1})}_{\tau_t} \sim p(\tau_t; \theta)} [r(s_t, a_t)]$$

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

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

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

$$E_{(x,y) \sim P(x,y)} [f(x)]$$

$$= \sum_x \sum_y f(x) p(x,y)$$

$$= \sum_x \sum_y f(x) p(y|x) p(x)$$

$$= \sum_x f(x) p(x) \underbrace{\sum_y p(y|x)}_1$$

$$= \sum_x f(x) p(x)$$

$$\mathbb{E}_{\tau \sim p(\tau; \theta)} \left[ \sum_{t=1}^T \left( \nabla_{\theta} \log p(a_t | s_t; \theta) \sum_{k=t}^T r(s_k, a_k) \right) \right]$$

$$= \mathbb{E} \left[ \begin{aligned} & \nabla_{\theta} \log p(a_1 | s_1; \theta) (r(s_1, a_1) + r(s_2, a_2)) \\ & + \nabla_{\theta} \log p(a_2 | s_2; \theta) r(s_2, a_2) \end{aligned} \right]$$

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

$$= \mathbb{E} \left[ \left( \nabla_{\theta} \log p(a_1 | s_1; \theta) + \nabla_{\theta} \log p(a_2 | s_2; \theta) \right) (r(s_1, a_1) + r(s_2, a_2)) \right]$$

$$= \mathbb{E} \left[ \begin{aligned} & \nabla_{\theta} \log p(a_1 | s_1; \theta) (r(s_1, a_1) + r(s_2, a_2)) \\ & + \nabla_{\theta} \log p(a_2 | s_2; \theta) (r(s_1, a_1) + r(s_2, a_2)) \end{aligned} \right]$$

$$\mathbb{E}_{\tau \sim p(\tau; \theta)} \left[ \nabla_{\theta} \log p(a_2 | s_2; \theta) r(s_1, a_1) \right] = 0$$

$$E_{\gamma \sim p(\gamma; \theta)} \left[ \nabla_{\theta} \log p(a_2 | s_2; \theta) \cdot r(s_1, a_1) \right] = ?$$

$$= \sum_{s_1, a_1} \sum_{s_2, a_2} \nabla_{\theta} \log p(a_2 | s_2; \theta) \cdot r(s_1, a_1) p(s_1, a_1, s_2, a_2; \theta)$$

$$= \sum_{s_1, a_1} r(s_1, a_1) p(s_1, a_1; \theta) \underbrace{\sum_{s_2, a_2} \nabla_{\theta} \log p(a_2 | s_2; \theta) p(s_2, a_2 | s_1, a_1; \theta)}$$

$$\sum_{s_2} \sum_{a_2} \nabla_{\theta} \log p(a_2 | s_2; \theta) p(a_2 | s_1, a_1, s_2; \theta) p(s_2 | s_1, a_1; \theta)$$

$$= \sum_{s_2} p(s_2 | s_1, a_1; \theta) \sum_{a_2} \nabla_{\theta} \log p(a_2 | s_2; \theta) p(a_2 | \underbrace{s_1, a_1, s_2}_{\text{redundant}}; \theta)$$

$$= \sum_{s_2} p(s_2 | s_1, a_1; \theta) \sum_{a_2} \nabla_{\theta} p(a_2 | s_2; \theta) \quad (\text{cond. ind.})$$

$$= \sum_{s_2} p(s_2 | s_1, a_1; \theta) \underbrace{\nabla_{\theta} \sum_{a_2} p(a_2 | s_2; \theta)}_1 = 0$$

## BASELINE SHIFT

variance of sampled  $\nabla_{\theta} J(\theta)$  is high, so we have to use many samples

subtract a baseline from  $r(r)$  to give an estimator with same mean and lower variance

$$\nabla_{\theta} J(\theta) = \nabla_{\theta} (J(\theta) - b)$$

$$\begin{aligned} &= \nabla_{\theta} E_{r \sim p(r; \theta)} [r(r) - b] \quad \leftarrow \text{since } E_{r \sim p(r; \theta)} [b] = b \\ &= E_{r \sim p(r; \theta)} \left[ \underbrace{\nabla_{\theta} \log p(r; \theta)}_{\text{call this } f(r)} (r(r) - b) \right] \end{aligned}$$

$$\begin{aligned} \text{Var}_{r \sim p(r; \theta)} [x] &= E_{r \sim p(r; \theta)} [x^2] - E_{r \sim p(r; \theta)} [x]^2 \\ &= E_{r \sim p(r; \theta)} \left[ f(r)^2 (r(r) - b)^2 \right] - (\nabla_{\theta} J(\theta))^2 \end{aligned}$$

↑  
this is bogus unless applied to each element  
of  $\nabla_{\theta}$  separately

## BASELINE SHIFT

$\Theta$  is a scalar

$$\begin{aligned}\text{var}_{\gamma \sim P(\gamma; \Theta)}[x] &= E_{\gamma \sim P(\gamma; \Theta)}[x^2] - E_{\gamma \sim P(\gamma; \Theta)}[x]^2 \\ &= E_{\gamma \sim P(\gamma; \Theta)} \left[ f(\gamma)^2 (r(\gamma) - b)^2 \right] - (\nabla_\Theta J(\Theta))^2\end{aligned}$$

minimize over  $b$ :

$$\frac{\partial}{\partial b} \text{var}[\cdot] = E_{\gamma \sim P(\gamma; \Theta)} \left[ -2 f(\gamma)^2 (r(\gamma) - b) \right] = 0$$

$$\rightarrow b = \frac{E[f(\gamma)^2 r(\gamma)]}{E[f(\gamma)^2]}$$

↑  
also look at  
2nd derivative ...

$$E[x x^T] = \text{cov}(x) + E[x] E[x]^T$$

$$E[f(\tau) (r(\tau) - b)]$$

$$\text{tr cov}[x] = E[\|x\|^2] - \|E[x]\|^2$$

$$= E[f(\tau)^T f(\tau) (r(\tau) - b)^2] - \|\nabla_{\theta} J(\theta)\|^2$$

$$\frac{\partial}{\partial b} (\text{tr cov}[x]) = E[-2 f(\tau)^T f(\tau) (r(\tau) - b)] = 0$$

$$\Rightarrow b = \frac{E[f(\tau)^T f(\tau) r(\tau)]}{E[f(\tau)^T f(\tau)]}$$



this at least makes sense  
but maybe isn't the best  
choice of baseline?

## IMPORTANCE SAMPLING (for off-policy learning)

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

$$= E_{x \sim q_b(x)} \left[ \left( \frac{p(x)}{q_b(x)} \right) f(x) \right]$$

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

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

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

## LIKELIHOOD RATIO

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