

598 RL

FALL 2020

POLICY GRADIENT
BRETL

POLICY GRADIENT - RUNNING EXAMPLE

$$S = \{0, 1\}$$

left right

$$A = \{0, 1\}$$

stay still move

with probability $\beta \in [0, 1]$
 we take the action that we
did not intend

rewards

if we end up in front of the left arm

+ 10



if $s_0 = 0$ and $a_0 = 0$

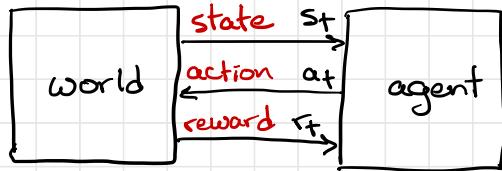
otherwise

+ 0

with prob $1 - \beta$, $s_1 = 0$ and $r_0 = 0$

with prob β , $s_1 = 1$ and $r_0 = 0$

MODEL



$p(s_0)$ = probability that the initial state is s_0

$p(s_{t+1} | s_t, a_t)$ = 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 θ

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

Tabular policy for finite state and action space

θ_{as} = weight of action a in state s

θ_{00}
↑
state 0
action 0

$\theta_{00}, \theta_{10}, \theta_{01}, \theta_{11}$

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

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

assume the policy is time-invariant

~~$p(a|s; \theta) = \theta_{as}$~~

s.t.

$$\theta_{as} \in [0, 1]$$

and

$$\sum_{a'} \theta_{a's} = 1$$

$$p(a=0 | s=0; \theta) = \frac{e^{\theta_{00}}}{e^{\theta_{00}} + e^{\theta_{10}}}$$

$$p(a=1 | s=0; \theta) = \frac{e^{\theta_{10}}}{e^{\theta_{00}} + e^{\theta_{10}}}$$

$$P(a=0|s=0; \theta) + P(a=1|s=0; \theta) = 1$$

$$\nabla_{\theta} () = \nabla_{\theta} (1) = 0$$

$$\nabla_{\theta} P(a=0|s=0; \theta) + \nabla_{\theta} P(a=1|s=0; \theta) = 0$$

GOAL

$$\tau = (s_0, a_0, \dots, s_{T-1}, a_{T-1}, s_T) \quad \leftarrow \text{trajectory}$$

$$p(\tau; \theta) = p(s_0) \prod_{t=0}^{T-1} p(a_t | s_t; \theta) p(s_{t+1} | s_t, a_t)$$

↑ probability of generating this trajectory with a given policy

PAYOUTOFF

$$J(\theta) = E_{\tau \sim p(\tau; \theta)} \left[\sum_{t=0}^{T-1} r(s_t, a_t) \right]$$

WE WANT TO
MAXIMIZE THIS

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

$$p(\tau; \theta) = p(s_0) p(a_0 | s_0; \theta) p(s_1 | s_0, a_0) p(a_1 | s_1; \theta)$$

$$\begin{aligned} \nabla_\theta p(\tau; \theta) &= p(s_0) \nabla_\theta p(a_0 | s_0; \theta) p(s_1 | s_0, a_0) p(a_1 | s_1; \theta) \\ &\quad + p(s_0) p(a_0 | s_0; \theta) p(s_1 | s_0, a_0) \nabla_\theta p(a_1 | s_1; \theta) \end{aligned}$$

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

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

↓
 $= \sum_{\tau} \frac{r(\tau)}{P(\tau; \theta)}$
 ↑ probability of generating this trajectory
 ↑ total reward for this trajectory
 ↑ sum over all possible trajectories

$$\nabla_{\theta} J(\theta) = \sum_{\tau} \frac{r(\tau)}{P(\tau; \theta)} \nabla_{\theta} P(\tau; \theta)$$

↑ how good is τ ?
 ↑ what change in θ would make
 the probability of τ increase
 the fastest?

$(s_0, a_0, s_1, a_1) \rightarrow$ 16 possible trajectories

PROBLEM #1

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

This does not have the form

$$\left[\sum_{\tau} f(\tau) p(\tau; \theta) \right] \leftarrow E_{\tau \sim p(\tau; \theta)} [f(\tau)]$$

so it can't be approximated by sampling τ with θ as

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

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

PROBLEM #2

$$\nabla_{\theta} J(\theta) = \sum_{\tau} r(\tau) \nabla_{\theta} P(\tau; \theta)$$

$$\begin{aligned} p(\tau; \theta) &= p(s_0, a_0, \dots, s_{T-1}, a_{T-1}, s_T; \theta) \\ &= p(s_0) \prod_{t=0}^{T-1} p(a_t | s_t; \theta) p(s_{t+1} | s_t, a_t) \end{aligned}$$

$$\begin{aligned} \nabla_{\theta} P(\tau; \theta) &= p(s_0) \sum_{t=0}^{T-1} \left(\left(\prod_{\substack{k=0 \\ k \neq t}}^{T-1} p(a_k | s_k; \theta) p(s_{k+1} | s_k, a_k) \right) \right. \\ &\quad \left. \cdot \nabla_{\theta} p(a_t | s_t; \theta) p(s_{t+1} | s_t, a_t) \right) \end{aligned}$$

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 ...

$$\begin{aligned}\log p(\tau; \theta) &= \log \left(p(s_0) \prod_{t=0}^{T-1} p(a_t | s_t; \theta) p(s_{t+1} | s_t, a_t) \right) \\ &= \log p(s_0) + \sum_{t=0}^{T-1} \left(\log p(a_t | s_t; \theta) + \log p(s_{t+1} | s_t, a_t) \right)\end{aligned}$$

then ...

$$\nabla_{\theta} \log p(\tau; \theta) = \sum_{t=0}^{T-1} \nabla_{\theta} \log p(a_t | s_t; \theta)$$


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 !

$$\underbrace{\nabla_{\theta} \log p(\tau; \theta)}_{\text{what we want to find}} = \frac{1}{p(\tau; \theta)} \underbrace{\nabla_{\theta} p(\tau; \theta)}_{\text{what we are supposed to find}}$$

↓

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

↓

$$\begin{aligned}\nabla_{\theta} J(\theta) &= \sum_{\tau} r(\tau) \nabla_{\theta} p(\tau; \theta) \\ &= \sum_{\tau} (r(\tau) \nabla_{\theta} \log p(\tau; \theta)) p(\tau; \theta) \\ &= E_{\substack{\tau \sim p(\tau; \theta)}} \left[r(\tau) \nabla_{\theta} \log p(\tau; \theta) \right]\end{aligned}$$

↑
can be computed by sampling τ with θ

THE PAYOFF GRADIENT

$$\nabla_{\theta} J(\theta) = E_{\tau \sim p(\tau; \theta)} \left[\left(\sum_{t=0}^{T-1} \nabla_{\theta} \log p(a_t | s_t; \theta) \right) \left(\sum_{t=0}^{T-1} r(s_t, a_t) \right) \right]$$

$$\approx \frac{1}{N} \sum_{i=1}^N \left(\sum_{t=0}^{T-1} \nabla_{\theta} \log p(a_t^i | s_t^i; \theta) \right) \left(\sum_{t=0}^{T-1} r(s_t^i, a_t^i) \right)$$

$$\tau^i = (s_0^i, a_0^i, \dots, s_{T-1}^i, a_{T-1}^i, s_T^i)$$

is generated with θ

$$l(\theta) = \frac{1}{N} \sum_{i=1}^N \sum_{t=0}^{T-1} \log p(a_t^i | s_t^i; \theta) r(\tau^i)$$

$$\approx \frac{1}{n} \sum_{i=1}^N \nabla_{\theta} \log p(\tau^i; \theta) r(\tau^i)$$

$$\approx \frac{1}{N} \sum_{i=1}^N \sum_{t=0}^{T-1} \nabla_{\theta} \log p(a_t^i | s_t^i; \theta) r(\tau^i)$$

REINFORCE algorithm

① sample τ^1, \dots, τ^N by running the policy $p(a_t | s_t; \theta)$

② estimate $\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \left(\sum_{t=0}^{T-1} \underbrace{\nabla_{\theta} \log p(a_t^i | s_t^i; \theta)}_{\text{red bracket}} \right) \left(\sum_{t=0}^{T-1} r(s_t^i, a_t^i) \right)$

③ update $\theta_{k+1} = \theta_k + \alpha \nabla_{\theta} J(\theta_k)$

repeat until convergence

↑
how to compute this?