

Policy gradient methods (PG methods)

The main idea behind PG methods is the "likelihood ratio" trick.

An illustration of this trick in a very simple setting.

Let  $X$  be a r.v. with mass function

$$p(\theta, \cdot)$$

$\xrightarrow{\text{parameter}}$

i.e.,  $p(X=x)$  is parameterized by  $\theta$ .

$$\mathcal{T}(\theta) = E f(x)$$

Goal:  $\min_{\theta} \mathcal{T}(\theta)$

Want to find  $\theta$  but  $\mathcal{T}$  using a gradient method

$$\theta_{t+1} = \theta_t - \beta_t \nabla \mathcal{T}(\theta_t)$$

Need: " $\nabla \mathcal{T}(\theta)$ "

Use "likelihood ratio" trick, which is shown below.

$$\mathcal{T}(\theta) = \sum_x f(x) p(\theta, x) \quad \leftarrow \text{LOTUS}$$

$$\nabla \mathcal{T}(\theta) = \nabla \left( \sum_x f(x) p(\theta, x) \right)$$

need a few conditions  
to justify interchange of  $\sum$  &  $\nabla$

$$= \sum_x f(x) \nabla p(\theta, x)$$

mild regularity  
conditions  $\rightarrow$  usually 1st  
one invoke "dominated"

Convergence theorem

$$\nabla J(\theta) = \sum_x f(x) \nabla p(\theta, x)$$

$$= \sum_x \left( f(x) \frac{\nabla p(\theta, x)}{p(\theta, x)} \right) p(\theta, x)$$

$$\nabla J(\theta) = E \left( f \frac{\nabla p}{p} \right) = E(f \nabla \log p)$$

To get an estimate of  $\nabla J(\theta)$ , I can sample from  $\underline{x}$

$$\frac{\nabla p(\theta, x)}{p(\theta, x)} \rightarrow \text{likelihood ratio.}$$

$$\begin{aligned} \nabla \log p(\theta) \\ = \frac{\nabla p(\theta)}{p(\theta)} \end{aligned}$$

Connecting likelihood ratio to RL:

Fix an SSP problem with state space  $\mathcal{X}$ , action space  $\mathcal{A}$ .

$\overline{\Pi}_{\text{det}}$  : class of admissible stationary deterministic policies

$\{ \pi : \pi : \mathcal{X} \rightarrow \mathcal{A} \text{ & it is timeinvariant} \}$

$\xrightarrow{x} \boxed{\text{unparameterized policy}} \rightarrow \pi(x) \rightarrow \text{an action } \pi : \boxed{\begin{matrix} \text{states} \\ \hline \text{actions} \end{matrix}}$

For PGI method, we consider stationary randomized policies, i.e.,

$$\pi : \mathcal{X} \rightarrow \Delta(\mathcal{A})$$

$\rightarrow$  set of all distributions over the actions.  
For simplicity, assume all actions available in all states.

e.g.  $\pi_{\theta}(x, a) = \frac{\exp(h(\theta, x, a))}{\sum_b \exp(h(\theta, x, b))}$  randomized policy

$$\pi_{\theta}(x) = [\pi_{\theta}(x, a), \forall a \in \mathcal{A}]$$

distribution over  $\mathcal{A}$ .

Simple example for  $h$ :  $h(\theta, x, a) = \theta^T \phi(x, a)$

$$\pi_{\theta}(x, a) = \frac{\exp(\theta^T \phi(x, a))}{\sum_b \exp(\theta^T \phi(x, b))}$$

} Boltzmann distribution aka Soft-max

Assumption A1: Policy  $\pi_{\theta}$  is a continuously differentiable function of  $\theta$   
 a.s.  $\nabla \log \pi_{\theta}$  exists.

Note! Every  $\pi_{\theta}$  is identified by its parameter  $\theta \in \mathbb{R}^d$

Goal:  $\min_{\theta \in \Theta} J_{\pi_{\theta}}(x^*)$

→ Find an approximately optimal policy in the class of parametrized policies

& we want to find the best parameter in a class

$$\{ \pi_{\theta} \mid \theta \in \Theta \}$$

e.g.  $\Theta \subset \mathbb{R}^d$  from start state

Want to find a  $\theta^* \in \arg \min_{\theta \in \Theta} J_{\pi_{\theta}}(x^*)$

## Lecture 35\*

Today: An expression for  $\nabla_{\theta} J_{\pi_{\theta}}(x^0)$

So that we can do

$$\theta_{t+1} = \theta_t - \beta_t \hat{\nabla}_{\theta} J(\theta_t)$$

$\downarrow$

estimate of  $\nabla_{\theta} J_{\pi_{\theta}}(x^0)$

*Stochastic gradient algorithm*

$\min_{\theta} f(\theta), \quad \theta_{t+1} = \theta_t - \beta_t f'(\theta_t)$

$\theta_t \rightarrow \theta^* \text{ as } t \rightarrow \infty, \quad f'(\theta^*) = 0$

With policy gradients, we would want to find a

$$\theta^* \text{ s.t. } \nabla_{\theta} J_{\theta^*}(x^0) = 0$$

Catch! We need to know  $\nabla_{\theta} J_{\theta^*}(x^0)$

usual RL settings, closed form of  $\nabla_{\theta} J$

is not available &

has to be estimated from samples".

For a deterministic policy:  $Q_{\pi}(x, a) = g(x, a) + \sum_{x'} p_{x,a}(x') J_{\pi}(x')$

$$J_{\pi}(x) = Q_{\pi}(x, \pi(x))$$

With randomized policies:  $J_{\pi}(x) = \sum_a \pi(x, a) Q_{\pi}(x, a)$

$$Q_{\pi}(x, a) = g(x, a) + \sum_{x'} P_{x,a}(x') J_{\pi}(x')$$

need to average over actions  
since  $\pi$  is the distribution

Policy gradient theorem: (for SS  $P_s$ )

Start with

$$J_{\pi_{\theta}}(x^o) = \sum_a \pi(x^o, a) Q_{\pi_{\theta}}(x^o, a)$$

$$\nabla_{\theta} J_{\pi_{\theta}}(x^o) = \nabla \left( \sum_a \pi_{\theta}(x^o, a) Q_{\pi_{\theta}}(x^o, a) \right)$$

$$= \sum_a \left( \nabla \pi_{\theta}(x^o, a) Q_{\pi_{\theta}}(x^o, a) + \pi_{\theta}(x^o, a) \nabla Q_{\pi_{\theta}}(x^o, a) \right)$$

$$= \sum_a \left( \nabla \pi_{\theta}(x^o, a) Q_{\pi_{\theta}}(x^o, a) + \pi_{\theta}(x^o, a) \nabla \left( g(x^o, a) + \sum_{x'} P_{x,a}(x') J_{\pi_{\theta}}(x') \right) \right)$$

↓  
does not depend on  $\theta$

$$= \sum_a \left[ \nabla \pi_{\theta}(x^o, a) Q_{\pi_{\theta}}(x^o, a) + \pi_{\theta}(x^o, a) \sum_{x'} P_{x,a}(x') \nabla J_{\pi_{\theta}}(x') \right]$$

$$= \sum_a \left[ \nabla \pi_{\theta}(x^o, a) Q_{\pi_{\theta}}(x^o, a) + \pi_{\theta}(x^o, a) \sum_{x'} P_{x,a}(x') \left( \sum_{a'} \nabla \pi_{\theta}(x', a') Q_{\pi_{\theta}}(x', a') + \pi(x', a') \sum_{x''} P_{x', a'}(x'') \nabla J_{\pi_{\theta}}(x'') \right) \right]$$

$$= \sum_{x \in S} \sum_{k=0}^{\infty} P(x_k=x | x^o, \pi_{\theta}) \sum_a \nabla \pi_{\theta}(x, a) Q_{\pi_{\theta}}(x, a)$$

↓  
 Prob of going from  $x^0$  to  $x$  in  $k$  steps  
 while following the policy  $\pi_\theta$

$$\nabla_{\theta} J_{\pi_\theta}(x^0) = \sum_{x \in \mathcal{X}} \sum_{k=0}^{\infty} P(x_k=x | x^0, \pi_\theta) \sum_a \nabla \pi_\theta(x, a) Q_{\pi_\theta}(x, a)$$

↑  
 Policy gradient theorem. → occupancy measure  
 (can be normalized!)

With some abuse of notation, the policy gradient fn is written as

$$\nabla_{\theta} J_{\pi_\theta}(x^0) = \sum_x \nabla \pi_\theta(x) \sum_a \nabla \pi_\theta(x, a) Q_{\pi_\theta}(x, a)$$

$$= E_{\pi} \left( \sum_a \nabla \pi_\theta(x, a) Q_{\pi_\theta}(x, a) \right)$$

$X$  is a r.v. governed by a distribution that uses  $\pi^\theta$

$\nabla \pi_\theta$  is available in closed form (we fixed the parameterization)

$Q_{\pi_\theta}$  → estimated using a sample path

PG update:  $\theta_{t+1} = \theta_t - \beta_t \sum_a \nabla \pi_{\theta_t}(x, a) \hat{Q}_{\pi_\theta}(x, a)$

$\hat{Q}_{\pi_\theta} \rightarrow \text{estimate of } Q_{\pi_\theta}$

$$\begin{aligned}\nabla_\theta J_{\pi_\theta}(x) &= E_{\pi_\theta} \left( \sum_a \nabla \pi(x, a) Q_{\pi_\theta}(x, a) \right) \\ &= E_{\pi_\theta} \left( \sum_a \pi(x, a) Q_{\pi_\theta}(x, a) \frac{\nabla \pi(x, a)}{\pi(x, a)} \right) \\ &= E_{\pi_\theta} \left( \sum_a \pi(x, a) Q_{\pi_\theta}(x, a) \nabla \log \pi(x, a) \right)\end{aligned}$$

Suppose  $A$  is r.v. chosen wif the distribution  $\pi$

$$\rightarrow = E_{\pi_\theta} \left( Q_{\pi_\theta}(x, A) \nabla \log \pi(x, A) \right)$$

both  $x, A$  are random

$X$  is chosen from a distribution  
p. occupancy measure

$A$  is chosen wif  $\pi(x, \cdot)$

$$\begin{aligned}\sum_{i=1}^{\infty} p(i) h(i) &\xrightarrow{\text{def}} \mathbb{E}_{p(X=i)}[h(X)] \\ &= E(h(X))\end{aligned}$$

Suppose  $E(\hat{Q}_\theta) = Q_\pi$ . Then,

$$\theta_{t+1} = \theta_t - \beta_t \hat{Q}_t \nabla \log \pi_{\theta_t}(x_t, A_t)$$

REINFORCE algorithm

A more naïve version:

Fix  $\pi_{\theta_t}$ , simulate an episode  $(x_0, a_0, \dots, \overset{\downarrow}{x_T})$  terminated

Collect a sample of the total cost from this

episode, call it  $\hat{Q}_t \rightarrow$  Check  $\hat{Q}_t$  is an unbiased estimate of  $Q_\pi$

$$\theta_{t+1} = \theta_t - \beta_t \hat{Q}_t \nabla \log \pi_{\theta_t}(x_0, a_0)$$

$\downarrow$

$$a_t \sim \pi(x_0, \cdot)$$

Aside! Want to estimate  $Q_\pi(x, a) = E \left( \sum g(x_t, a_t) \middle| \begin{array}{l} x_0=x \\ a_0=a \end{array} \right)$

total cost in SSP

starting in  $x$  & taking action  $a$

Estimation!  
 $(x, a)$   
 $(x_0, a_0)$   
Follow  $\pi$  until termination.

Collect total cost sample, say  $C_t$

$$E(C_t) = Q_\pi(x, a)$$

Can be extended to cover an action  $A$  that is

chosen from a random policy  $\pi(\cdot, \cdot)$

Remark: In REINFORCE, policy evaluation ( $= Q^{\pi_0}(\cdot, \cdot)$ )  
 ↓ "Monte Carlo".

Instead, if we use a parametric approximation  
 for  $Q_\pi$ , say  $Q_{\pi_\theta}(x, a) \underset{\text{↑}}{\approx} \gamma^T \phi(x, a)$   
 linear func. approx.

Can we TD with LFA to obtain an estimate  
 of  $Q_{\pi_\theta}(\cdot, \cdot)$ .

Suppose TD converges to  $\gamma^*$ .

$$Q_\pi(x, a) \underset{\text{↑}}{\approx} \gamma^{*T} \phi(x, a)$$

$$E(\gamma^{*T} \phi(x, a)) \neq Q_\pi(x, a) ??$$

"Parametric approximation induce a bias"

On the other hand, MC methods usually have large variance.

Can do policy gradient with function approximation  
 ↓  
 Actor-Critic algorithms.

$$\theta_{t+1} = \theta_t - \beta_t \hat{\nabla}_{\theta_t} J_{\theta_t}(r^*)$$

↓

$$\hat{Q}_{\pi_0}(x_0, a_0) \cdot \nabla \log \pi(x_0, a_0)$$

↓

$$= r^{*\top} \phi(x_0, a_0)$$

How to get  $r^*$ ?

Fix policy  $\pi_{\theta_t}$

Obtain a sample path using  $\pi_{\theta_t}$

Run TD using this sample path  $\rightarrow \hat{r} \rightarrow$  estimate of  $r^*$

Use  $\hat{r}$  & do gradient descent in policy parameter.

