

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 $p(\theta, \cdot)$

\rightarrow parameter

i.e., $P(X=x)$ is parameterized by θ .

$$J(\theta) = E f(X)$$

e.g. $\theta \in \mathbb{R}^d$

Goal:

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

\rightarrow min among a class of parameterized r.v.s

Want to find best θ using a gradient method

$$\theta_{t+1} = \theta_t - \beta_t \nabla J(\theta_t) \quad \leftarrow \text{Gradient descent}$$

Need: " $\nabla J(\theta)$ "

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

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

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

need a few conditions
regularity
to justify interchange
of \sum & ∇

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

need regularity
conditions \rightarrow usually let
one r.v. be "dominated"

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

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

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

To get an estimate of $\nabla J(\theta)$, I can sample from \rightarrow

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

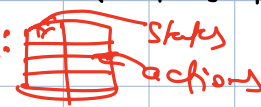
$\nabla \log p(\theta) = \frac{\nabla p(\theta)}{p(\theta)}$

Connecting likelihood ratio to RL:

Fix an SSP problem with state space X , action space A .

Π_{det} : class of admissible stationary deterministic policies

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

$\pi \rightarrow$ unparameterized policy $\rightarrow \pi(x) \rightarrow$ an action π : 

for Ph method, we consider stationary randomized policies, i.e.,

$\pi : X \rightarrow \Delta(A)$ \rightarrow set of all distributions over the actions.
 For simplicity, assume all actions available in all states.

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

is the probability of choosing action a in state x

randomized policy parameterized by θ

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

distribution over \mathcal{A}

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

$\theta \in \mathbb{R}^d$ *$\phi \in \mathbb{R}^d$*

state-action features

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

Boltzmann distribution aka Soft-max

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

Note: Every π_{θ} is identified by its parameter $\theta \in \mathbb{R}^d$

Goal: $\min_{\theta \in \Theta} J_{\pi_{\theta}}(x^0)$

Find an approximately optimal policy in the class of parameterized policies

& we want to find the best parameter in a class

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

e.g. $\Theta \subset \mathbb{R}^d$ *action state*

Want to find a $\theta^* \in \arg \min_{\theta \in \Theta} J_{\pi_{\theta}}(x^0)$

" J_{θ} " is not necessarily a convex function of parameter θ

So, a gradient also in policy space converges for local optima.

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

So that we can do

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

← Policy gradient update iteration

↓
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 not to find a θ^* s.t. $\nabla_{\theta^*} J(x^0) = 0$

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

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

need to average over actions since π is distribution

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

Policy gradient theorem: (for SSPs)

Start with

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

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

$$= \sum_a \left(\nabla_{\theta} \pi_{\theta}(x^0,a) Q_{\pi_{\theta}}(x^0,a) + \pi_{\theta}(x^0,a) \nabla_{\theta} Q_{\pi_{\theta}}(x^0,a) \right) \quad \text{--- (1)}$$

$$= \sum_a \left(\nabla_{\theta} \pi_{\theta}(x^0,a) Q_{\pi_{\theta}}(x^0,a) + \pi_{\theta}(x^0,a) \nabla_{\theta} \left(g(x^0,a) + \sum_{x'} P_{x,a}(x') J_{\pi_{\theta}}(x') \right) \right)$$

does not depend on θ

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

$$= \sum_a \left[\nabla_{\theta} \pi_{\theta}(x^0,a) Q_{\pi_{\theta}}(x^0,a) + \pi_{\theta}(x^0,a) \sum_{x'} P_{x,a}(x') \nabla_{\theta} J_{\pi_{\theta}}(x') \right]$$

$$= \sum_a \left(\nabla_{\theta} \pi_{\theta}(x^0,a) Q_{\pi_{\theta}}(x^0,a) + \right.$$

$x_0 \rightarrow$ pos. prob $\pi_{\theta}(x^0,a)$
 $x' \rightarrow$ pos. prob $P_{x^0,a}(x')$

$$\pi_{\theta}(x^0,a) \sum_{x'} P_{x^0,a}(x') \left(\sum_{a'} \nabla_{\theta} \pi_{\theta}(x',a') Q_{\pi_{\theta}}(x',a') + \pi_{\theta}(x',a') \sum_{x''} P_{x',a'}(x'') \nabla_{\theta} J_{\pi_{\theta}}(x'') \right)$$

$$= \sum_{x \neq x^0} \sum_{k=0}^{\infty} P(x_k = x | x^0, \pi_{\theta}) \sum_a \nabla_{\theta} \pi_{\theta}(x,a) Q_{\pi_{\theta}}(x,a)$$

↓
 Prob of going from x^0 to x in k steps
 while following the policy π_θ

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

↑
 Policy gradient theorem.

→ occupancy measure
 (can be normalized)

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

$$\nabla_{\theta} J_{\pi_{\theta}}(x^0) = \sum_x \underbrace{d^{\pi_{\theta}}(x)}_{= \sum_{k=0}^{\infty} P(x_k = x | x^0, \pi_{\theta})} \sum_a \nabla \pi_{\theta}(x, a) Q_{\pi_{\theta}}(x, a)$$

$$= \mathbb{E}_{\pi} \left(\sum_a \nabla \pi_{\theta}(X, a) Q_{\pi_{\theta}}(X, a) \right)$$

↓
 X is a r.v. governed by a distribution that uses $d^{\pi_{\theta}}$

$\nabla \pi_{\theta}$ is available in closed form (we had the parametrization)

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

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

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

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

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

Suppose A is a r.v. (chosen using the distribution π)

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

both X, A are random

X is chosen from a distribution using the occupancy measure

A is chosen using $\pi(X, \cdot)$

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

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

REINFORCE algorithm

$$\begin{aligned}\sum_{i=1}^{100} p(i) h(i) &\rightarrow \text{inf } p(x=i) \\ &= E(h(X))\end{aligned}$$

A more naïve version:

Fix π_{θ_t} , Simulate an episode $(x_0, a_0, \dots, \overset{\text{terminal}}{\downarrow} x_T)$

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_{π}

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

\downarrow
 $A_0 \sim \pi(x_0, \cdot)$

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

total cost in SSP

Starting in x & taking action a

Estimation! $(\overset{x}{\parallel} \overset{a}{\parallel} \overset{\text{red highlight}}{x_0, a_0, \dots})$

\downarrow
follow π 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 $\pi(x, \cdot)$

Remark: In REINFORCE, policy evaluation (^{Estimating} $= Q^{\pi_\theta}(\cdot, \cdot)$)
is "Monte Carlo".

Instead, if we use a parametric approximation
for Q_π , say $Q_{\pi_\theta}(x, a) \approx r^\top \phi(x, a)$
 \uparrow
Linear func. approx.

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

Suppose TD converges to r^* .

$$Q_\pi(x, a) \approx r^{*\top} \phi(x, a)$$

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

"Parametric approximation induce a bias"

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

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

$$\theta_{t+1} = \theta_t - \beta_t \underbrace{\hat{\nabla}_{\theta} J_{\theta_t}(x^0)}_{\downarrow} \\ \underbrace{\hat{Q}_{\pi_0}(x_0, A_0) \cdot \nabla \log \pi(x_0, A_0)}_{\downarrow} \\ = \gamma^{*\top} \phi(x_0, A_0)$$

How to get γ^* ?

Fix policy π_{θ_t}

Obtain a sample path using π_{θ_t}

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

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

