

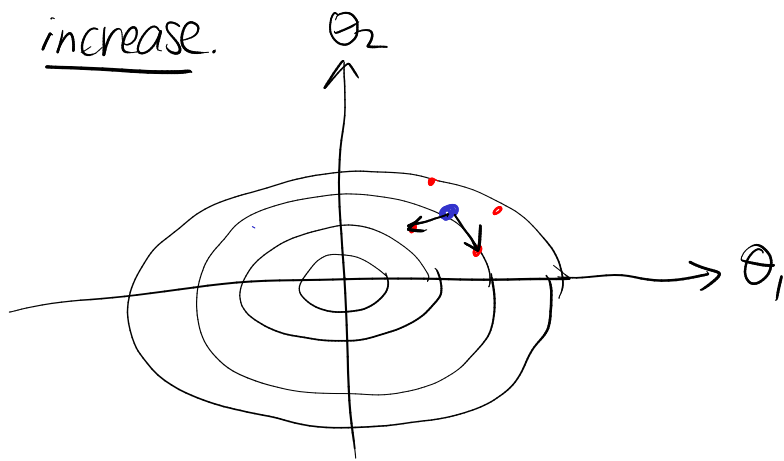
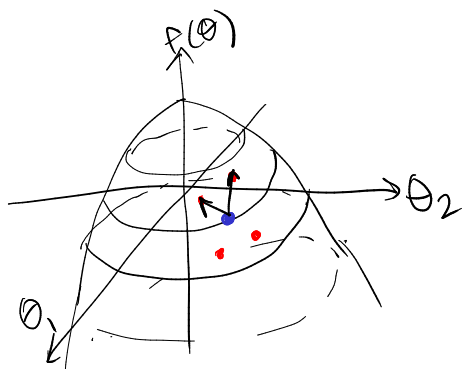
# Lecture 14: Policy Gradients

## 1) Derivative-Free Optimization

How can we find maxima only using function evaluation?  
 i.e. we can query  $f(\theta): \mathbb{R}^d \rightarrow \mathbb{R}$  but not  $\nabla f(\theta)$ .

Goal: find a descent direction

Simple idea: randomly test a few directions & see  
 which lead to increase.



There are many variations of this simple idea:  
 simulated annealing, cross-entropy method, genetic algorithms,  
 evolutionary strategies. They differ in how random samples are  
 aggregated into update step.

We will cover methods that use samples to construct gradient estimates.

### A) Random Search

Recall when we discussed iQR the finite difference approximation:

$$f'(x) \approx \frac{f(x+\delta) - f(x-\delta)}{2\delta}$$

This idea can help us build an approximation of the  
gradient based only on function evaluation.  
 ↙ direction of steepest ascent

For vector functions:  $\theta \in \mathbb{R}^d$

$$\langle \nabla f(\theta), v \rangle \approx \frac{f(\theta + \delta v) - f(\theta - \delta v)}{2\delta}$$

Alg: Random Search

initialize  $\theta_0$

for  $t=0, 1, \dots$

sample  $v_1, \dots, v_N \sim N(0, \mathbb{I})$

update  $\theta_{t+1} = \theta_t + \frac{\alpha}{2\delta N} \sum_{k=1}^N (f(\theta_t + \delta v_k) - f(\theta_t - \delta v_k)) v_k$

We can understand this as stochastic gradient ascent:

$$\begin{aligned} \mathbb{E}((f(\theta + \delta v_k) - f(\theta - \delta v_k)) / (2\delta)) &\approx \mathbb{E}(2\delta \nabla f(\theta)^T v_k \cdot v_k) \\ &= 2\delta \mathbb{E}[v_k v_k^T] \nabla f(\theta) \\ &= 2\delta \nabla f(\theta) \end{aligned}$$

This method samples/searches in parameter space. ( $\theta$ )

### B) Importance Weighting

Distribution trick: in general, we can write:

$$f(\theta) = \mathbb{E}_{x \sim P_\theta} [h(x)]$$

for some class of distributions  $P_\theta$

(In RL setting,  $P_\theta$  could represent the distribution over trajectories induced by  $\pi_\theta$ .)

Now suppose a sampling distribution  $\rho$  where  $\frac{P_\theta(x)}{\rho(x)} < \infty$ .

$$\mathbb{E}_{x \sim P_\theta} [h(x)] = \sum_x h(x) P_\theta(x) \cdot \frac{\rho(x)}{\rho(x)} = \mathbb{E}_{x \sim \rho} \left[ \frac{P_\theta(x)}{\rho(x)} h(x) \right]$$

↑  
"importance weights"

This allows us to write the gradient:

$$\nabla f(\theta) = \mathbb{E}_{x \sim p} \left[ \frac{\nabla_{\theta} P_{\theta}(x)}{p(x)} h(x) \right]$$

This is true for any  $p(x)$ . If we pick  $p(x) = P_{\theta}(x)$  then

$$\nabla_{\theta} f(\theta) = \mathbb{E}_{x \sim P_{\theta}(x)} \left[ \frac{\nabla_{\theta} P_{\theta}(x)}{P_{\theta}(x)} h(x) \right] = \mathbb{E}_{x \sim P_{\theta}(x)} \left[ \nabla_{\theta} \log(P_{\theta}(x)) h(x) \right]$$

Now if  $P_{\theta}(x)$  factors,  $\log(P_{\theta}(x))$  will be sum of factors, and the gradient will depend only on factors which depend on optimization variable (This is very useful for policy optimization)

Therefore, our stochastic maximization algorithm:

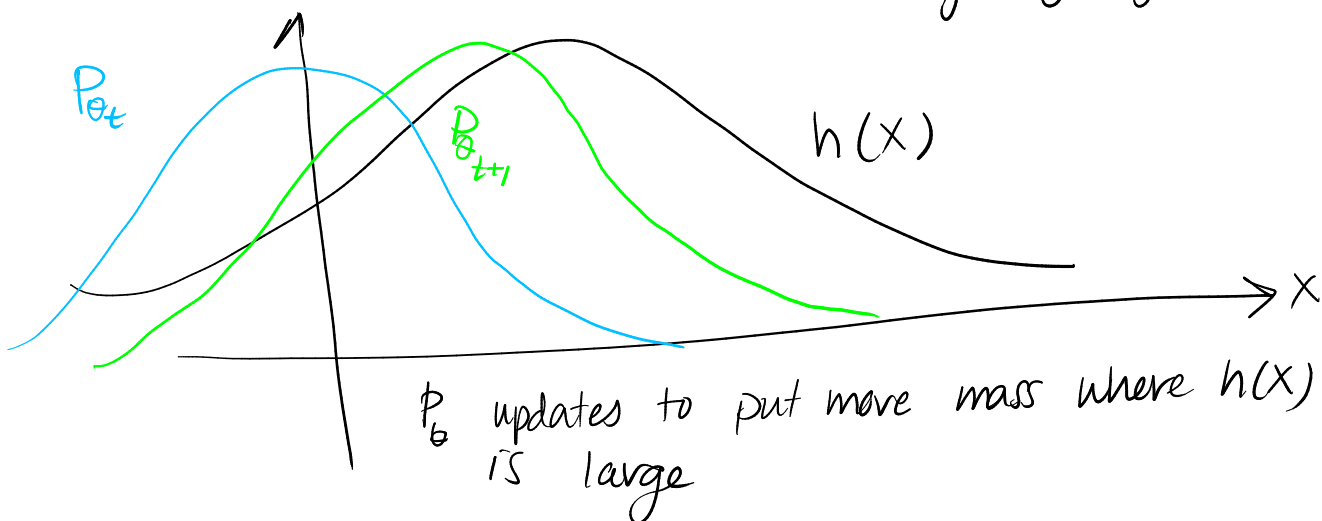
Alg: sampling-DFD  
initialize  $\theta_0$

for  $t=0, 1, \dots$   
sample  $x_i \sim P_{\theta_t}$  and observe  $h(x)$   $i=1, \dots, N$

$$\theta_{t+1} = \theta_t + \frac{\alpha}{N} \sum_{i=1}^N \nabla_{\theta_t} \log(P_{\theta_t}(x_i)) h(x_i)$$

This method samples in  $x$ -space rather than parameter space.

↖ e.g. trajectory space



## 2) Policy Optimization via Simple Random Search

Recall the RL setting:

MDP  $\mathcal{M} = \{S, A, P, r, \gamma\}$  with  $P, r$  unknown

parametrized policy  $\pi_\theta$ ,  $\theta \in \mathbb{R}^d$  (e.g. weights of neural networks)

objective function:  $J(\theta) = \mathbb{E}_{s_0 \sim \mu_0} \left[ \sum_{t=0}^{\infty} \gamma^t r_t \mid P, r, \pi_\theta \right]$

observe "rollout" of  $\pi_\theta$ : trajectory  $\tau = (s_0, a_0, s_1, \dots)$  and rewards  $(r_0, r_1, \dots)$

means we can "sample"  $J(\theta) = \mathbb{E}_{\tau} (R(\tau))$  & observe  $\tau, R(\tau)$

### Meta-Algorithm: Derivative-Free SGA

initialize  $\theta_0$

for  $t=0, 1, \dots$

1) collect rollouts using  $\theta_t$

2) compute estimate  $g_t$  of  $\nabla_{\theta} J(\theta)$  using rollouts

3)  $\theta_{t+1} = \theta_t + \alpha g_t$

The rest of lecture: 3 ways to estimate  $\nabla J(\theta)$  using rollouts.

### Simple Random Search

Based on the "random search" idea.

1) collect rollouts:  $\tau^+$  &  $\tau^-$  with

$\pi_{\theta_t + \delta v}$  &  $\pi_{\theta_t - \delta v}$  for  $v \sim \mathcal{N}(0, 1)$ , small  $\delta > 0$

2) compute estimate:  $g_t = \frac{1}{2\delta} (R(\tau^+) - R(\tau^-)) v$

### 3) Policy Gradient (PG) from Trajectories (REINFORCE)

Another approach based on "importance weighting" derivation.

$$\tau = (s_0, a_0, s_1, \dots) \text{ and } p_\theta(\tau) = \mu_\theta(s_0) \pi_\theta(a_0|s_0) P(s_1|s_0, a_0) \pi_\theta(a_1|s_1) \dots$$

$$J(\theta) = \mathbb{E}_{\tau \sim p_\theta} [R(\tau)]$$

Claim: for  $\tau \sim p_\theta(\tau)$  (i.e.  $\tau$  observed from rollout with  $\pi_\theta$ )

$g = \sum_{t=0}^{\infty} \nabla_{\theta} [\log(\pi_{\theta}(a_t|s_t))] R(\tau)$  is an unbiased estimate of  $\nabla J(\theta)$ .

Proof: Using derivation from earlier in lecture with  $f \leftarrow J$ ,  $h \leftarrow R$ ,  $x \leftarrow \tau$ :

$$\nabla J(\theta) = \mathbb{E}_{\tau \sim p_\theta} [\nabla_{\theta} [\log(p_\theta(\tau))] R(\tau)]$$

$$\begin{aligned} \nabla_{\theta} \log(p_\theta(\tau)) &= \nabla_{\theta} [\log(\mu_\theta(s_0) \pi_\theta(a_0|s_0) P(s_1|s_0, a_0) \pi_\theta(a_1|s_1) \dots)] \\ &= \nabla_{\theta} [\cancel{\log(\mu_\theta(s_0))} + \sum_{t=0}^{\infty} \log(\pi_\theta(a_t|s_t)) + \cancel{\log(P(s_{t+1}|s_t, a_t))}] \\ &= \sum_{t=0}^{\infty} \nabla_{\theta} \log(\pi_\theta(a_t|s_t)) \quad \text{no } \theta \text{ dependence} \quad \square \end{aligned}$$

$\nabla_{\theta} \log p_\theta(\tau)$  ends up not depending at all on unknown transition function  $P$ !

REINFORCE: 1) collect rollout  $\tau$  with  $\pi_\theta$

2) compute estimate

$$g = \sum_{t=0}^{\infty} \nabla_{\theta} \log(\pi_\theta(a_t|s_t)) R(\tau)$$

# 4) Policy Gradient with value functions

PG w/ trajectories often has high variance. An alternative commonly used in practice uses an alternative estimate using Q functions.

Claim: for  $s, a \sim d_{\pi_\theta}^{s_0}$ ,

$$g = \frac{1}{1-\gamma} \nabla_\theta \log(\pi_\theta(a|s)) Q^{\pi_\theta}(s, a)$$

is an unbiased estimate of  $\nabla J(\theta)$

Proof:  $\nabla J(\theta) = \nabla_\theta \mathbb{E}_{s_0 \sim \mathcal{M}} [V^{\pi_\theta}(s_0)]$

value fn. def.

$$= \mathbb{E}_{s_0 \sim \mathcal{M}} [\nabla_\theta \mathbb{E}_{a_0 \sim \pi_\theta(s_0)} [Q^{\pi_\theta}(s_0, a_0)]]$$

$s_0$  doesn't depend on  $\theta$  & Q fn. def.

$$\begin{aligned} \nabla_\theta \mathbb{E}_{a_0 \sim \pi_\theta(s_0)} [Q^{\pi_\theta}(s_0, a_0)] &= \sum_{a_0 \in \mathcal{A}} \nabla_\theta [\pi(a_0|s_0) Q^{\pi_\theta}(s_0, a_0)] \quad \text{defn. of expectation} \\ &= \sum_{a_0 \in \mathcal{A}} (\nabla_\theta \pi(a_0|s_0)) Q^{\pi_\theta}(s_0, a_0) + \pi(a_0|s_0) \nabla_\theta Q^{\pi_\theta}(s_0, a_0) \\ &\quad \left\{ \begin{array}{l} \text{importance} \\ \text{weighting} \\ \text{trick} \end{array} \right. \quad \left\{ \begin{array}{l} r(s_0, a_0) \text{ doesn't} \\ \text{depend on } \theta \end{array} \right. \\ &= \mathbb{E}_{a_0 \sim \pi_\theta(s_0)} [\nabla_\theta \log \pi_\theta(a_0|s_0) Q^{\pi_\theta}(s_0, a_0)] + \gamma \mathbb{E}_{\substack{a_0 \sim \pi_\theta(s_0) \\ s_1 \sim P(s_0, a_0)}} [\nabla_\theta V^{\pi_\theta}(s_1)] \end{aligned}$$

$$\nabla J(\theta) = \mathbb{E}_{\substack{s_0 \sim \mathcal{M}_0 \\ a_0 \sim \pi_\theta(s_0)}} [\nabla_\theta \log \pi_\theta(a_0|s_0) Q^{\pi_\theta}(s_0, a_0)] + \gamma \mathbb{E}_{s_1 \sim P_1^{\pi_\theta}} [\nabla_\theta V^{\pi_\theta}(s_1)]$$

$$\mathbb{E}_{s_0 \sim \mathcal{M}_0} [\nabla_\theta V^{\pi_\theta}(s_0)]$$

we can iterate!

$$\nabla J(\theta) = \sum_{t=0}^{\infty} \gamma^t \mathbb{E}_{s_t, a_t \sim \pi_{\theta}^t} \left[ \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \cdot Q^{\pi_{\theta}}(s_t, a_t) \right]$$

expanding  
expectation

$$= \sum_{t=0}^{\infty} \sum_{\substack{s \in \mathcal{S} \\ a \in \mathcal{A}}} P_t^{\pi_{\theta}}(s, a; \mu_0) \gamma^t \cdot \nabla_{\theta} \log \pi_{\theta}(a | s) \cdot Q^{\pi_{\theta}}(s, a)$$

definition  
of  $d_{\mu_0}^{\pi_{\theta}}$

$$= \frac{1}{1-\gamma} \mathbb{E}_{s, a \sim d_{\mu_0}^{\pi_{\theta}}} \left[ \nabla_{\theta} \log \pi_{\theta}(a | s) \cdot Q^{\pi_{\theta}}(s, a) \right]$$

□

One final gradient estimate:  $s, a \sim d_{\mu_0}^{\pi_{\theta}}$   
 $\frac{1}{1-\gamma} \nabla_{\theta} \log \pi_{\theta}(a | s) \cdot (Q^{\pi_{\theta}}(s, a) - b(s))$

Baseline function  $b(s)$  further helps in variance reduction. Most common  $b(s) = V^{\pi_{\theta}}(s)$  results in advantage function-based PG

$$A^{\pi_{\theta}}(s, a) = Q^{\pi_{\theta}}(s, a) - V^{\pi_{\theta}}(s).$$

Exercise: show that  $\mathbb{E}_{a \sim \pi_{\theta}(s)} \left[ \nabla_{\theta} \log \pi_{\theta}(a | s) \cdot b(s) \right] = 0$

for any action-independent baseline.