

# Lecture 22: Reinforcement Learning I

Hallmark of RL is a collection of algorithms

for "black-box optimization" of stochastic optimal control problems

→  $\sum g(x,u) \dots$   
only have  
access to the costs

don't have  $f(x,u) \dots$

$$\sum_1^\infty g(x,u) \rightarrow \sum_{n=0}^\infty E[g(x,u)]$$

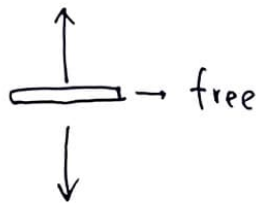
$$E\left[\sum_{n=0}^\infty g(x,u)\right]$$

What we optimize.

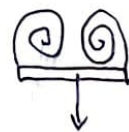
How do you optimize?

Example. Fluid dynamics of flapping flight

The heaving foil (Jun Zhang, NYU)



symmetric flat plate



unstable flow



dynamics governed  
by P.D.E.

difficult to compute dynamics

The Policy Gradient "trick" (REINFORCE)

$$\min_{\alpha} E[g(x)] \quad \text{with} \quad x \sim p_{\alpha}(x)$$

take gradient

$$\frac{d}{d\alpha} E[g(x)] = \frac{d}{d\alpha} \int dx g(x) P_{\alpha}(x)$$

$$= \int dx g(x) \frac{d}{d\alpha} P_{\alpha}(x)$$

$$\frac{d}{d\alpha} \frac{d}{d\alpha} \log P_{\alpha}(x) = \frac{1}{P_{\alpha}(x)} \frac{d}{d\alpha} P_{\alpha}(x)$$

$$= \int dx g(x) \frac{d}{d\alpha} \log P_{\alpha}(x) \cdot P_{\alpha}(x)$$

$$= E \left[ g(x) \frac{d}{d\alpha} \log P_{\alpha}(x) \right]$$

approx with Monte-carlo

$$\approx \frac{1}{N} \sum_{n=1}^N g(x_n) \frac{d}{d\alpha} \log P_{\alpha}(x_n)$$

in optimal control case,

$$\frac{d}{d\alpha} E \left[ \sum_{n=0}^N g(x[n], u[n]) \right] = \sum_{n=0}^N \int dx[n] du[n]$$

$$g(x[n], u[n]) \frac{d}{d\alpha} P_{\alpha}(x[n], u[n])$$

$$= E \left[ \sum_{n=0}^N g(x[n], u[n]) \frac{d}{d\alpha} \log P_{\alpha}(x[n], u[n]) \right]$$

$$P_{\alpha}(x[n], u[n]) = P_0(x[n]) \prod_{k=1}^n P(x[k] | x[k-1], u[k-1])$$

$$\prod_{k=0}^n P_{\alpha}(u[k] | x[k])$$

Taking the log we have

$$\log P_{\alpha}(x[n], u[n]) = \log P_0(x[n]) + \sum_{k=1}^n \log P(x[k] | x[k-1], u[k-1])$$

$$+ \sum_{k=0}^n \log P_{\alpha}(u[k] | x[k])$$

Only the last term depends on  $\alpha$ , which yields

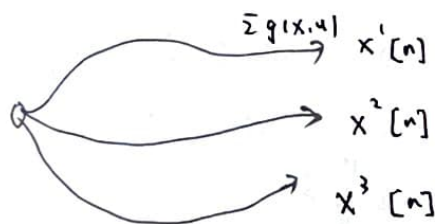
$$\frac{d}{d\alpha} E \left[ \sum_{n=0}^N g(x[n], u[n]) \right] = E \left[ \sum_{n=0}^N g(x[n], u[n]) \sum_{k=0}^n \frac{d}{d\alpha} \log P_{\alpha}(u[k] | x[k]) \right]$$

gradient on controller, not dependent on  $P(x[k] | u[k-1], x[k-1])$   
surprise?

but a weak statement, because optimize expert?

Intuition

$N$  random trajectories



$$u = -Kx + w$$

↑  
random number

make trajectory with  
better (lower) cost  
higher probability

Black-box optimization

$$\min_{\alpha} g(\alpha)$$

How do you do gradient-free optimization?

Idea: finite differences

$$\frac{\partial g}{\partial \alpha} \bigg|_{\alpha} \approx \begin{bmatrix} \vdots & \vdots & \vdots & \vdots & \vdots \\ g(\alpha + \epsilon_1) - g(\alpha) & \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots & \vdots \end{bmatrix}$$

↑  
n parameters

↑  
n+1 evaluations

expensive evaluation  
index  $L$

Then gradient descent:  $\alpha[k+1] = \alpha[k] - \eta \frac{\partial g}{\partial \alpha}^T \bigg|_{\alpha[k]}$

## Stochastic gradient descent

key idea: as long as I'm going "downhill", can get away w/ less evaluations of  $g(\alpha)$

Simpler idea: "Weight perturbation"

$$\Delta \alpha = -\eta [g(\alpha + \beta) - g(\alpha)] \beta$$

small random vector  
 $B \sim \mathcal{N}(0, \sigma)$

2 evaluations

$g(\alpha + \beta) > g(\alpha)$ , in decrease direction

$g(\alpha + \beta) < g(\alpha)$ , the same decrease direction

$$E[\Delta \alpha] = -\eta \sigma^2 \frac{\partial g}{\partial \alpha}^T$$

$$\Delta \alpha = -\frac{\eta}{\sigma^2} [g(\alpha + \beta) - b] \beta$$

↑  
baseline

"expected performance"

still

$$E[\Delta \alpha] \propto -\frac{\partial g}{\partial \alpha}^T$$

↑  
proportional

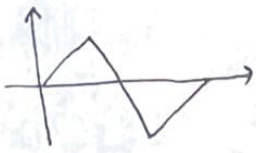
Key limitation of RL today "Sample Complexity"

Policy parameters for the heaving foil



Fourier base, not a good idea  
different changes to the  
cost.

Cost function: max inverse cost of transport



Converged reliably to a triangle wave

very repeatable experiments

"extremum-seeking control"

"iterative learning control"