

# 20. Reinforcement Learning 2

## Toxonomies of RL

- Methods based on value functions
  - $V^\pi(x)$  is expected long term reward of starting at  $x$  and executing policy  $\pi$
  - Policy Search methods; Actor-critic Methods
  - Policy-Gradient method  $\subset$  Policy Search
  - PPO, is Actor-critic method
  - Policy gradient as optimizer
- Topic 2
- Topic 3

## Today's Topic outline

- Reinforce  $\subset$  Policy Gradient
  - $\min_{\alpha} E[f(x)], x \propto p_{\alpha}(x)$ , distribution of  $x$ ,  $N(\alpha, \sigma(\text{fixed}))$
  - write a distribution over possible  $x$ 's and minimize the expected value
  - delta-like function is the goal, highest possibility at lowest value
  - optimize of gradient descent:  $\frac{\partial}{\partial \alpha} E[f(x)] = E[f(x) \frac{\partial}{\partial \alpha} \ln P_{\alpha}(x)]$ , log-likelihood method (or policy gradient "trick")
- Reinforce with additive Gaussian noise
  - $x \sim P_{\alpha} \sim N(\alpha, \sigma^2)$
  - $x = \alpha + \beta, \beta \sim N(0, \sigma^2)$
  - $P_{\alpha}(x) = C e^{-\frac{(x-\alpha)^T(x-\alpha)}{2\sigma^2}}$ , probability density function of gaussian
  - $\ln P_{\alpha}(x) = \frac{-(x-\alpha)^T(x-\alpha)}{2\sigma^2} + \dots$  terms that do not depend on  $\alpha$
  - $\frac{\partial}{\partial \alpha} \ln P_{\alpha}(x) = \frac{1}{\sigma^2}(\alpha - x)^T = \frac{1}{\sigma^2}\beta^T$
  - $f(x) \frac{\partial}{\partial \alpha} \ln P_{\alpha}(x) = \frac{1}{\sigma^2} f(\alpha + \beta) \beta^T$
  - $\Delta \alpha = -\eta \frac{1}{\sigma^2} f(\alpha + \beta) \beta^T$ ,  $\eta$ : learning rate

- given a small perturbation:
- $\Delta\alpha = -\eta \frac{1}{\sigma^2} [f(\alpha + \beta) - f(\alpha)] \beta^T$ 
  - if  $f(\alpha + \beta) > f(\alpha)$ , move  $-\beta$  direction
  - if  $f(\alpha + \beta) < f(\alpha)$ , move  $+\beta$  direction
- If you have gradients, why not use them? (from AutoDiff)
  - the answer is subtle!
  - scienrio: a wsg gripper try to grip a brick,  $z_{height}$ 
    - controller: descend until  $z_{close}$ , close gripper, raise hand
    - rewards: height of brick at time = 5 sec
    - plot reward in Y vs  $z_{close}$  in X, very discontinuous loss landscapes
    - gradient descent on discontinuous landscapes in general doesn't work very well
    - **but adding probability density function, the smoothing effect works well**
- The idea of Non-smooth optimization
  - is "randomized smoothing"
  - new interpretation that, policy gradient in RL is sort of 1 to 1 mapping Randomized Smoothing
  - example,  $\min |x|_1$ , l1 norm
- In RL the randomrization comes from
  - $P_\alpha(x)$ , exploration
  - Random initial conditions
  - Domain randomrization
  - the Smoothing effect helps convergence and optimization

## 1. Good papers on RL

- Do Differentiable Simulators Give Better Policy Gradients?

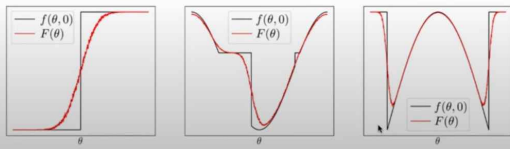


[Do Differentiable Simulators Give Better Policy Gradients.pdf](#)

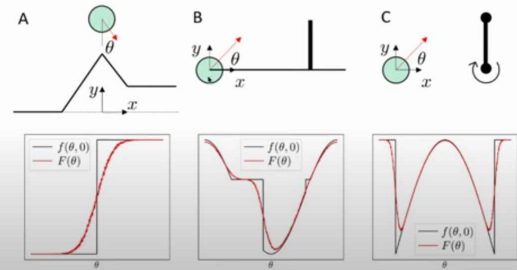
### Smoothing with stochasticity

$$\min_{\theta} f(\theta) \quad \text{vs} \quad \min_{\theta} E_w [f(\theta, w)]$$

$$w \sim N(0, \Sigma)$$



### Smoothing with stochasticity for Multibody Contact



- Do Differentiable Simulators give better policy gradients? the answer is subtle

### Randomized smoothing

- Approximate smoothed objective via Monte-carlo :

$$E_{\mu} [f(x)] \approx \frac{1}{K} \sum_{i=1}^K f(x_i), \quad x_i \sim \mathcal{N}(\mu, \Sigma)$$

- First-order gradient estimate

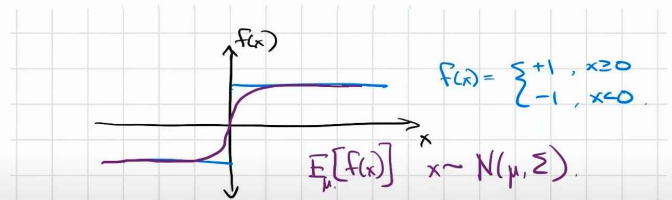
$$\frac{\partial}{\partial \mu} E_{\mu} [f(x)] \approx \frac{1}{K} \sum_{i=1}^K \frac{\partial f(\mu + w_i)}{\partial \mu}, \quad w_i \sim \mathcal{N}(0, \Sigma)$$

J. Burke, F. E. Curtis, A. Lewis, M. Overton, and L. Simoes, *Gradient Sampling Methods for Nonsmooth Optimization*, 02 2020, pp. 201–225.

- Zero-order gradient estimate (aka REINFORCE)

$$\frac{\partial}{\partial \mu} E_{\mu} [f(x)] \approx \frac{1}{K} \sum_{i=1}^K [f(\mu + w_i) - f(\mu)] w_i, \quad w_i \sim \mathcal{N}(0, \Sigma)$$

### Example: The Heaviside function



First-order estimator is biased

$$\frac{\partial f(x)}{\partial x} = 0 \text{ almost everywhere!}$$

- Is

## 2. What does policy gradient look like for control? (open question)

- linear Gaussian dynamics + quadratic cost

## 3. Examples