

Lecture 23: Reinforcement Learning II

RL is black-box optimization for stochastic optimal control

Last time Policy Gradient methods

$$u = \pi_{\alpha}(x)$$

↑

parameter vector

e.g. $u = -Kx$

$$u = \text{DNN}(x)$$

$$u = \pi_{\alpha}(y)$$

"output feedback"

↑ outputs

$$\dot{x} = f(x, u), \quad y = h(x, u)$$

Trial-and-error

Too many concerns:

1) sample complexity

2) Non-convexity

$$\dot{x} = Ax + Bu$$

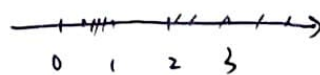
$$y = Cx$$

$$u = -Ky$$

$$A = \begin{bmatrix} 0 & 0 & 2 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix}, \quad B = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$$

$$C = [1 \quad 1 \quad 3]$$

stable K



set of good controllers
are disconnected

tasks diverse enough that bad local minimum
disappeared, domain randomization

Today - Q-Learning

- Actor-Critic

Idea: Learning Value Functions

Policy Gradient

* Dynamic Programming

← inefficient → efficient

evaluate your dynamics
exactly once from the
terminal back to $t=0$

Learn $f(x,u)$	Learn $\pi(x)$	Learn $\hat{J}(x)$
-------------------	-------------------	-----------------------

system i.d.

Policy
Search

cost-to-go

+
model-based
control

$J^*(x), J^\pi(x)$

$$J(x) = E \left[\sum_{n=0}^{\infty} \gamma^n g(x[n], u[n]) \right] \quad x[n+1] = f(x[n], u[n], w[n])$$

\uparrow int. horizon \nwarrow discount factor $0 \leq \gamma \leq 1$

Learn only $\hat{J}^*(x)$

$$J^*(x) = \min_u g(x,u) + \gamma E \left[\hat{J}^*(f(x,u)) \right]$$

depends on knowing $f(x,u)$

$$\pi^*(x) = \arg \min_u [\quad]$$

Two work-arounds

- 1) Q-learning $Q(x, u)$ instead of $J(x)$
- 2) Actor-critic Learn $\pi(x)$ and $J^*(x)$ simultaneously

Q function

Define $Q^\pi(x, u) = E[g(x, u) + J^\pi(f(x, u))]$

$$J^*(x) = \min_u Q^\pi(x, u)$$

J has $\dim(x)$ inputs (scalar ^{output} ~~input~~)

Q has $\dim(x) + \dim(u)$ inputs (scalar output)

$$\pi^*(x) = \arg \min_u Q^*(x, u)$$

Q is deep neural network DNN

QT-Opt says runtime optimization of $\hat{Q}^*(x, u)$

$\left(\begin{array}{l} Q^*(x, u) \text{ can be difficult to model} \\ \rightarrow \text{google arm farm} \end{array} \right.$

Actor-Critic

Policy gradient

$$\Delta \alpha = -\gamma [J^{\pi^x}(\alpha + \beta) - b] \left[\frac{\partial}{\partial \alpha} \log f_{\pi^x}(x|\alpha) \right]^T$$

one time experiment \nearrow

Idea: use $J^{\pi^x}(\alpha)$ \nwarrow policy gradient from REINFORCE

reestimate $\pi(x)$
 \rightarrow update $J^\pi(x)$

as the baseline

Expected value
(better baseline than running an extra trail)

Toddler

feedback, linear comb. nonlinear basis

dead beat controller, return every cycle

batch update does not change mean, but reduce variance

if reset policy, keep value function, converge to
policy fast

Online learning only

