

Lecture 7: Policy gradient

7.1 Motivation

$$\text{previous} \left\{ \begin{array}{l} V_{\bar{a}}(s) \sim V_{\theta}(s) \\ Q_{\bar{a}}(s, a) \sim Q_{\theta}(s, a) \\ \text{policy} \sim \underline{\epsilon\text{-greedy}} \end{array} \right.$$

In this lecture.

$$\pi(s, a) \leftarrow \pi_{\theta}(s, a) = \text{prob}(a | s, \theta)$$

parameterize policy.

Advantage

- David Silver
- ① Better convergence
 - ② Effective in high-dim / continuous action space
 - ③ can learn stochastic properties

- Sutton
- ④ a good way to inject prior knowledge.
 - <1> 线性
 - <2> 非线性

Disadvantage.

- ① Typically converge to local minima
- ② Evaluating a policy \rightarrow often \rightarrow inefficient, high variance.

Ex 1. rock - paper - scissors
 deterministic \rightarrow exploit \rightarrow fail
 random (Nash equilibrium)

Ex 2.

\leftarrow		\downarrow		\leftarrow
X		¥		X

① shadow area $\pi_{\theta} \in \mathcal{R}^n$
 these two grids use same policy.

7.2. policy objective functions

Goal: Find π_{θ} (θ), find best θ

Measure quality of π_{θ}

① start value.

$$J_v(\theta) = V^{\pi_{\theta}}(s_1) = E_{\pi_{\theta}}(V_1)$$

② average value.

$$J_{av} V(\theta) = \sum_s d(s) V^{\pi_{\theta}}(s)$$

③ average reward per time-step.

$$J_{av} R(\theta) = \sum_s d(s) \sum_a \pi_{\theta}(s, a) R_s^a$$

$d(s)$ 是在 π_{θ} 下 s 的稳态分布.

7.2.1 Ex: one-step MDP

(prob. of (s, a) is 0.25)

$$\begin{aligned}
 \nabla_{\theta} J(\theta) &= \sum_s d(s) \sum_a \nabla_{\theta} \pi_{\theta}(s, a) R_s^a \\
 &= \sum_s d(s) \sum_a \pi_{\theta}(s, a) \frac{\nabla_{\theta} \pi_{\theta}(s, a)}{\pi_{\theta}(s, a)} R_s^a \\
 &= \sum_s d(s) \sum_a \pi_{\theta}(s, a) \nabla_{\theta} \ln \pi_{\theta}(s, a) R_s^a \\
 &= E_{s, a} [\nabla_{\theta} \ln \pi_{\theta}(s, a) R_s^a] \\
 &= E_{\pi_{\theta}} [\nabla_{\theta} \ln \pi_{\theta}(s, a) R_s^a] \\
 \Delta \theta &= \underline{\nabla_{\theta} \ln \pi_{\theta}(s, a) R_s^a}
 \end{aligned}$$

Score function: $\nabla_{\theta} \ln \pi_{\theta}(s, a)$

Ex2. $\pi_{\theta}(s, a) \propto e^{\phi(s, a)^T \theta}$

$$\pi_{\theta}(s, a) = \frac{1}{\sum_a e^{\phi(s, a)^T \theta}} e^{\phi(s, a)^T \theta} \quad \text{prob}(a|s, \theta)$$

$$\begin{aligned}
 \nabla_{\theta} \ln \pi_{\theta}(s, a) &= \nabla_{\theta} \ln \left(\frac{1}{\sum_a e^{\phi(s, a)^T \theta}} e^{\phi(s, a)^T \theta} \right) \\
 &= \nabla_{\theta} \ln e^{\phi(s, a)^T \theta} - \frac{1}{\sum_a e^{\phi(s, a)^T \theta}} \sum_a e^{\phi(s, a)^T \theta} \phi(s, a) \\
 &= \phi - \langle \phi \rangle
 \end{aligned}$$

7.3 REINFORCE

$$J(\theta) = V_{\pi_{\theta}}(s)$$

key is to find $\nabla J(\theta)$

Policy Gradient Theorem

For $J_v, J_{avR}, \frac{1}{1-\gamma} J_{avV}$

$$\nabla_{\theta} J(\theta) = E_{\pi_{\theta}} \left(\underbrace{\nabla_{\theta} \ln \pi_{\theta}(s, a)}_{\text{policy}} \cdot \underbrace{Q^{\pi_{\theta}}(s, a)}_{\text{value}} \right)$$

Prove: $\gamma=1$

$$\nabla V_{\pi}(s) = \nabla \left(\sum_a \pi(a|s) q_{\pi}(s, a) \right)$$

$$= \sum_a \nabla \pi(a|s) \cdot q_{\pi}(s, a) + \pi(a|s) \cdot \nabla q_{\pi}(s, a)$$

$$= \sum_a \nabla \pi(a|s) \cdot q_{\pi}(s, a) + \pi(a|s) \nabla \left(\sum_{s', r} P(s', r|s, a) (r + V_{\pi}(s')) \right)$$

$$= \sum_a \nabla \pi(a|s) \cdot q_{\pi}(s, a) + \pi(a|s) \nabla \left(\sum_{s'} P(s'|s, a) V_{\pi}(s') \right)$$

$$= \sum_a \nabla \pi(a|s) \cdot q_{\pi}(s, a) + \underbrace{\pi(a|s) \sum_{s'} P(s'|s, a)}_{\sum_{s'} \sum_a \pi(a|s) P(s'|s, a) \dots} \nabla V_{\pi}(s')$$

$$= \sum_a \nabla \pi(a|s) \cdot q_{\pi}(s, a) + \pi(a|s) \sum_{s'} P(s'|s, a)$$

$$\left[\sum_{a'} \nabla \pi(a'|s') \cdot q_{\pi}(s', a') + \pi(a|s) \sum_{s''} P(s''|s', a') \nabla V_{\pi}(s'') \right]$$

⋮

$$= \sum_{x \in S} \left(\sum_{k=0}^{\infty} P(S \rightarrow x, k, \pi) \right) \sum_a \nabla \bar{u}(a|x) \frac{\partial \bar{u}(x, a)}{\partial a}$$

$$= \sum_x \underline{d(x)} \sum_a \nabla \bar{u}(a|x) \frac{\partial \bar{u}(x, a)}{\partial a}$$

$$\nabla J(\theta) = \nabla V_{\bar{u}}(s_0)$$

$$= \sum_s \underline{d(s)} \sum_a \nabla \bar{u}(a|s) \frac{\partial \bar{u}(s, a)}{\partial a}$$

$$= \sum_s \underline{d(s)} \sum_a \pi(a|s) \nabla (\ln \bar{u}(a|s)) \cdot \frac{\partial \bar{u}(s, a)}{\partial a}$$

$$= \bar{E}_{S, a} \left(\nabla \ln \bar{u}(a|s) \cdot \frac{\partial \bar{u}(s, a)}{\partial a} \right)$$

$$= \bar{E}_{\bar{u}} \left(\nabla \ln \bar{u}(a|s) \cdot \frac{\partial \bar{u}(s, a)}{\partial a} \right)$$

注意!!!

Monte-Carlo Policy Gradient (REINFORCE)

- Update parameters by stochastic gradient ascent
- Using policy gradient theorem
- Using return v_t as an unbiased sample of $Q^{\pi_\theta}(s_t, a_t)$

$$\Delta \theta_t = \alpha \nabla_{\theta} \log \pi_{\theta}(s_t, a_t) v_t$$

function REINFORCE

Initialise θ arbitrarily

for each episode $\{s_1, a_1, r_2, \dots, s_{T-1}, a_{T-1}, r_T\} \sim \pi_{\theta}$ **do**

for $t = 1$ to $T - 1$ **do**

$\theta \leftarrow \theta + \alpha \nabla_{\theta} \log \pi_{\theta}(s_t, a_t) v_t$

end for

end for

return θ

end function

REINFORCE: Monte-Carlo Policy-Gradient Control (episodic) for π_*

Input: a differentiable policy parameterization $\pi(a|s, \theta)$

Algorithm parameter: step size $\alpha > 0$

Initialize policy parameter $\theta \in \mathbb{R}^{d'}$ (e.g., to $\mathbf{0}$)

Loop forever (for each episode):

Generate an episode $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$, following $\pi(\cdot|\cdot, \theta)$

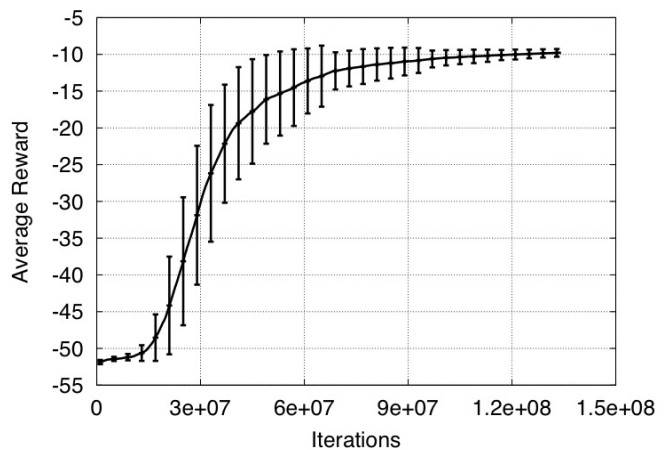
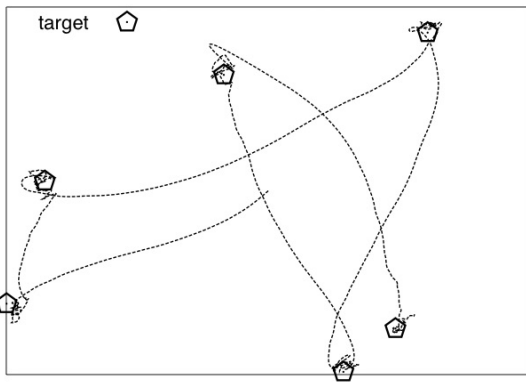
Loop for each step of the episode $t = 0, 1, \dots, T - 1$:

$$G \leftarrow \sum_{k=t+1}^T \gamma^{k-t-1} R_k \quad (G_t)$$

$$\theta \leftarrow \theta + \alpha \gamma^t G \nabla \ln \pi(A_t | S_t, \theta)$$

$\mu \neq 1$

Puck World Example



- Continuous actions exert small force on puck
- Puck is rewarded for getting close to target
- Target location is reset every 30 seconds
- Policy is trained using variant of Monte-Carlo policy gradient

☆: MC has high variance.

7.4 Reduce Variance.

① use a Base line.

$$\nabla J(\theta) \propto \sum_s d(s) \sum_a \left(\frac{q}{\pi} (s, a) - b(s) \right) \nabla \pi(a|s, \theta)$$

$$\left[\sum_s d(s) \sum_a \frac{q}{\pi} (s, a) \nabla \bar{a}(a|s, \theta) \right] - \underbrace{\sum_s d(s) \sum_a b(s) \nabla \pi(a|s, \theta)}_{= 0}$$

$$\begin{aligned} & \sum_s d(s) b(s) \sum_a \nabla \bar{a}(a|s, \theta) \\ = & \sum_s d(s) b(s) \nabla \sum_a \pi(a|s, \theta) \\ = & 0 \end{aligned}$$

REINFORCE with Baseline (episodic), for estimating $\pi_\theta \approx \pi_*$

Input: a differentiable policy parameterization $\pi(a|s, \theta)$

Input: a differentiable state-value function parameterization $\hat{v}(s, \mathbf{w})$

Algorithm parameters: step sizes $\alpha^\theta > 0$, $\alpha^\mathbf{w} > 0$

Initialize policy parameter $\theta \in \mathbb{R}^d$ and state-value weights $\mathbf{w} \in \mathbb{R}^d$ (e.g., to $\mathbf{0}$)

Loop forever (for each episode):

Generate an episode $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$, following $\pi(\cdot|\cdot, \theta)$

Loop for each step of the episode $t = 0, 1, \dots, T-1$:

$$G \leftarrow \sum_{k=t+1}^T \gamma^{k-t-1} R_k \quad (G_t)$$

$$\delta \leftarrow G - \hat{v}(S_t, \mathbf{w})$$

$$\mathbf{w} \leftarrow \mathbf{w} + \alpha^\mathbf{w} \delta \nabla_{\hat{v}} \hat{v}(S_t, \mathbf{w})$$

$$\theta \leftarrow \theta + \alpha^\theta \gamma^t \delta \nabla_{\pi} \ln \pi(A_t|S_t, \theta)$$

②. Using Critic

Critic = update action-value function

actor = update policy parameters

Action-Value Actor-Critic

- Simple actor-critic algorithm based on action-value critic
- Using linear value fn approx. $Q_w(s, a) = \phi(s, a)^\top w$
 - Critic** Updates w by linear TD(0)
 - Actor** Updates θ by policy gradient

function QAC

Initialise s, θ

Sample $a \sim \pi_\theta$

for each step **do**

Sample reward $r = \mathcal{R}_s^a$; sample transition $s' \sim \mathcal{P}_s^a$.

Sample action $a' \sim \pi_\theta(s', a')$

$\delta = r + \gamma Q_w(s', a') - Q_w(s, a)$

$\theta = \theta + \alpha \nabla_\theta \log \pi_\theta(s, a) Q_w(s, a)$

$w \leftarrow w + \beta \delta \phi(s, a)$

$a \leftarrow a', s \leftarrow s'$

end for

end function

- The **policy gradient** has many equivalent forms

$\nabla_{\theta} J(\theta) = \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(s, a) v_t]$	REINFORCE
$= \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(s, a) Q^w(s, a)]$	Q Actor-Critic
$= \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(s, a) A^w(s, a)]$	Advantage Actor-Critic
$= \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(s, a) \delta]$	TD Actor-Critic
$= \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(s, a) \delta e]$	TD(λ) Actor-Critic

$$A^w(s, a) = Q^{\bar{w}_{\theta}}(s, a) - V^{\bar{w}_{\theta}}(s, a)$$

$$\delta = R + \gamma V(s', w) - V(s, w)$$