

Reinforcement Learning: Policy Gradient

AI/ML Teaching

Goals & Keyword

- Policy-Based RL vs Value-Based RL
- Monte-Carlo Policy Gradient
- Actor-Critic Policy Gradient

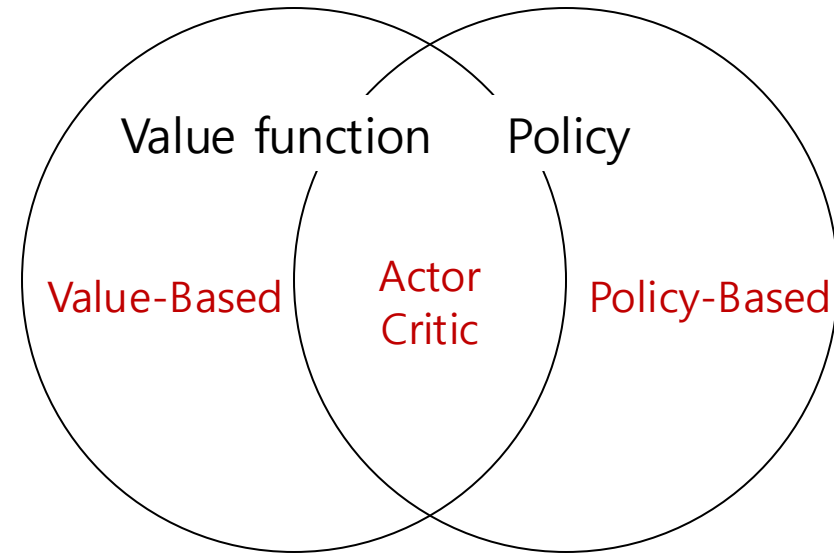
Policy-Based RL vs Value-Based RL

- Value-based RL: A policy is generated directly from the value function, $V^\pi(s)$ or $Q^\pi(s, a)$
 - e.g. ϵ -greedy
- Policy-based RL: Directly parameterize the policy

$$\pi(s, a) = \mathbb{P}[a|s]$$

Policy-Based RL vs Value-Based RL

- Value Based
 - Learnt value function
 - Implicit policy
- Policy Based
 - No value function
 - Learnt policy
- Actor-Critic
 - Learnt value function
 - Learnt policy



Policy-Based RL vs Value-Based RL

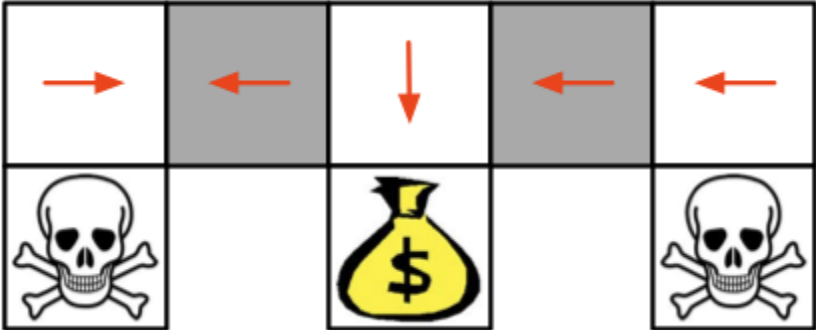
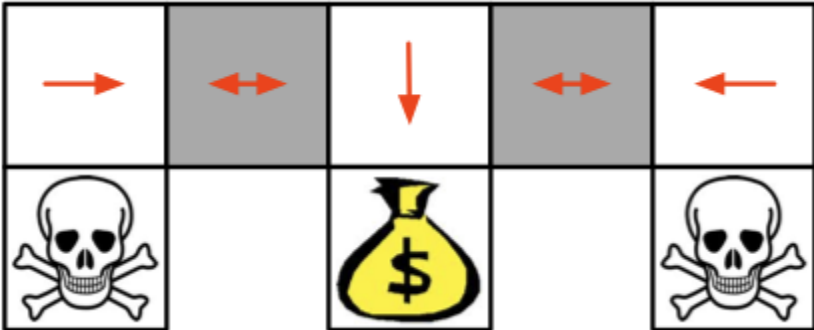
- Policy-Based RL
 - (+) Better convergence properties
 - (+) Effective in high-dimensional or **continuous** action spaces
 - (+) Can learn **stochastic** policies
 - (-) typically converge to a local rather than global optimum*
 - (-) Evaluating a policy is typically inefficient and high variance

*Value-based RL could also converge to local optimum (function approx., bootstrapping, off-policy)

Stochastic policy

- Stochastic environment
- Non-fully observable state



Deterministic policy	Stochastic policy
<p>*Agent cannot differentiate the grey states</p>  <p>A 2x5 grid representing a deterministic policy. The top row contains five cells with red arrows: right, left, down, left, left. The bottom row contains three cells with icons: a skull and crossbones, a yellow money bag with a dollar sign, and a skull and crossbones. The second and fourth cells in both rows are shaded grey.</p>	 <p>A 2x5 grid representing a stochastic policy. The top row contains five cells with red arrows: right, double-headed horizontal, down, double-headed horizontal, left. The bottom row contains three cells with icons: a skull and crossbones, a yellow money bag with a dollar sign, and a skull and crossbones. The second and fourth cells in both rows are shaded grey.</p>

Policy Search

Policy Objective Functions

- In episodic environments we can use the start value

$$J_1(\theta) = V^{\pi_\theta}(s_1) = \mathbb{E}_{\pi_\theta}[v_1]$$

- In continuing environments, we can use the average value

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

- Or the average reward per time-step

$$J_{avR}(\theta) = \sum_s d^{\pi_\theta}(s) \sum_a \pi_\theta(s, a) \mathcal{R}_s^a$$

where $d^{\pi_\theta}(s)$: stationary distribution of Markov chain for π_θ

Policy gradient

- Find θ maximizing $J(\theta)$: policy ascent
- Numerical policy gradient $\frac{\partial J(\theta)}{\partial \theta_k} \approx \frac{J(\theta + \epsilon u_k) - J(\theta)}{\epsilon}$
- Compute Policy gradient analytically: Policy gradient theorem

- We can compute $\nabla_{\theta} \pi_{\theta}(s, a)$

$$\nabla_{\theta} \pi_{\theta}(s, a) = \pi_{\theta}(s, a) \frac{\nabla_{\theta} \pi_{\theta}(s, a)}{\pi_{\theta}(s, a)} = \pi_{\theta}(s, a) \nabla_{\theta} \log \pi_{\theta}(s, a)$$

(def) Score function

Policy gradient theorem

$$\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(s, a) Q_{\pi}(s, a)]$$

(shortened derivation)

$$\begin{aligned} \nabla_{\theta} J(\pi_{\theta}) &= \nabla_{\theta} \sum_{s \in \mathcal{S}} d(s) V^{\pi}(s) \\ &= \sum_{s \in \mathcal{S}} d(s) \nabla_{\theta} V^{\pi}(s) \\ &= \sum_{s \in \mathcal{S}} d(s) \nabla_{\theta} \left(\sum_{a \in \mathcal{A}} Q^{\pi}(s, a) \pi_{\theta}(s, a) \right) \\ &= \sum_{s \in \mathcal{S}} d(s) \sum_{a \in \mathcal{A}} Q^{\pi}(s, a) \nabla_{\theta} \pi_{\theta}(s, a) \\ &= \sum_{s \in \mathcal{S}} d(s) \sum_{a \in \mathcal{A}} Q^{\pi}(s, a) \pi_{\theta}(s, a) \nabla_{\theta} \log \pi_{\theta}(s, a) \\ &= \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(a|s) Q_{\pi}(s, a)] \end{aligned}$$

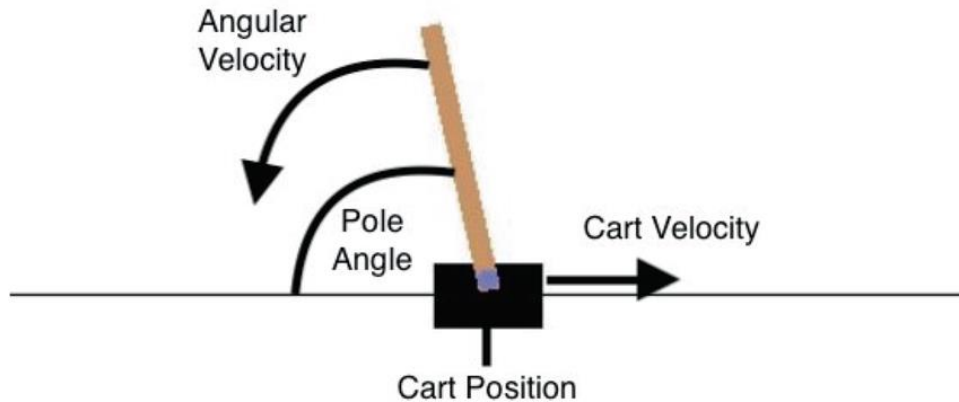
(Example) Softmax policy

- Cartpole
 - Action: left/right
 - Softmax policy

$$\text{Left: } \pi_{\theta}(a_0|s) = \frac{\exp(h_{\theta}(a_0|s))}{\sum_{a \in \mathcal{A}} \exp(h_{\theta}(a|s))} = p$$

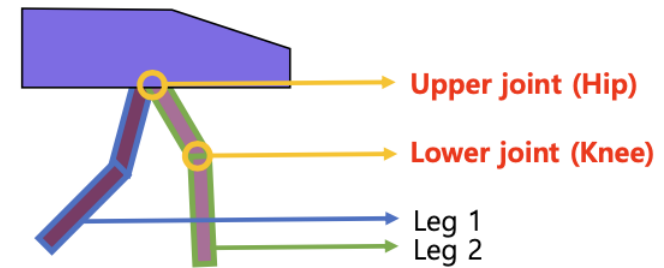
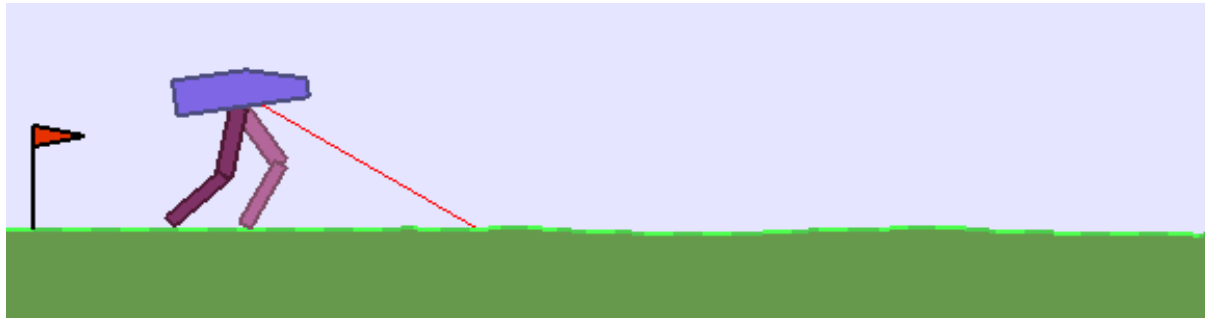
$$\text{Right: } \pi_{\theta}(a_1|s) = \frac{\exp(h_{\theta}(a_1|s))}{\sum_{a \in \mathcal{A}} \exp(h_{\theta}(a|s))} = 1 - p$$

Action: sampling $(p/1 - p) \rightarrow$ left or right



(Example) Gaussian policy

- Bipedal walker
 - Action: motor speed values for each of the 4 joints at both hips and knees



$$\pi_{\theta}(a|s) = \frac{1}{\sqrt{2\pi}\sigma(s, \theta)} \exp\left(-\frac{(a - \mu(s, \theta))^2}{2\sigma(s, \theta)^2}\right)$$

Joint1: $\mu_{\theta_1}, \sigma_{\theta_1}$ from h_{θ_1}
Joint2: $\mu_{\theta_2}, \sigma_{\theta_2}$ from h_{θ_2}
Joint3: $\mu_{\theta_3}, \sigma_{\theta_3}$ from h_{θ_3}
Joint4: $\mu_{\theta_4}, \sigma_{\theta_4}$ from h_{θ_4}

Joint1: sample from $\mathcal{N}(\mu_{\theta_1}, \sigma_{\theta_1})$
Joint2: sample from $\mathcal{N}(\mu_{\theta_2}, \sigma_{\theta_2})$
Joint3: sample from $\mathcal{N}(\mu_{\theta_3}, \sigma_{\theta_3})$
Joint4: sample from $\mathcal{N}(\mu_{\theta_4}, \sigma_{\theta_4})$

Monte-Carlo Policy Gradient (REINFORCE)

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

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

function REINFORCE

 Initialize θ

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

Actor-Critic

- Monte-Carlo policy gradient has high variance
- We use a critic to estimate the action-value function
- Actor-critic algorithms maintain two sets of parameters
 - Critic: updates action-value function parameters w
 - Actor: updates policy parameters θ
- Actor-critic algorithms follow an approximate policy gradient

$$\begin{aligned}\nabla_{\theta} J(\theta) &\approx \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(s, a) Q_w(s, a)] \\ \Delta \theta &= \alpha \nabla_{\theta} \log \pi_{\theta}(s, a) Q_w(s, a)\end{aligned}$$

Actor-Critic

function QAC

Initialize 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 $a' \sim \pi_\theta(s', a')$

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

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

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

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

end for

Reference

- David Silver, COMPM050/COMPGI13 Lecture Notes
- Richard S. Sutton and Andrew G. Barto, "Reinforcement Learning: An Introduction," 2nd Ed.
- 김재훈, "Introduction to Policy Gradient," DMQA Seminar