

Policy Gradient

- Basic 분류 Policy를 구분하기



* Policy Based RL Find θ that maximises $J(\theta)$

장점) better convergence properties

effective in high-dimensional action spaces

can learn stochastic policies

단점) converge to local

high variance

$$\theta_{t+1} = \theta_t + \Delta \theta$$

$$= \theta_t + \alpha \nabla_{\theta} J(\theta)$$

$$= \theta_t + \alpha \left(\mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(s, a) Q^{\pi_{\theta}}(s, a)] \right) \quad \begin{matrix} \nearrow \text{Policy Gradient} \\ \text{Theorem} \end{matrix}$$

- Monte-Carlo Policy Gradient (REINFORCE) (θ).

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

using stochastic gradient ascent, policy gradient theorem

- Actor-Critic Policy Gradient (θ, w)

< Actor (update action-value function)

Critic (to estimate action-value function)

update

parameter θ

parameter w

$$\Delta \theta = \alpha \nabla_{\theta} \log \pi_{\theta}(s, a) Q_w(s, a) \quad (\text{approximate policy gradient})$$

- Policy Gradient algorithms $\frac{\partial J}{\partial \theta}$ find $\theta \leftarrow$ use $J(\theta)$!

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

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

$$\mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(s, a) \delta_c]$$

TD(λ) Actor-Critic

$$G_{\theta}^{-1} \nabla_{\theta} J(\theta) = \omega$$

Natural Actor-Critic