CS-866 Deep Reinforcement Learning

Improved Policy-Based Learning IV: Deterministic Policy Gradient



Nazar Khan
Department of Computer Science
University of the Punjab

Deterministic Policy Gradient (DPG)

Actor-Critic recap:

- ► Actor-critic methods combine policy and value learning
- Can reduce variance and improve learning performance

Idea behind deterministic policy gradient:

- Learn a value function $Q_{\phi}(s,a)$
- Use it as a differentiable target to optimize the policy
- ► Policy *follows* the value function

Reference:1

¹David Silver et al. 'Deterministic policy gradient algorithms'. In: *International Conference on Machine Learning*. 2014, pp. 387–395.

Deterministic Policy Gradient Objective

Collect data D and train a value network $Q_{\phi}(s, a)$.

Then optimize a deterministic policy $\pi_{\theta}(s)$ via:

$$J(heta) = \mathbb{E}_{s \sim D} \left[\sum_{t=0}^n Q_\phi(s, \pi_ heta(s))
ight]$$

Chain rule gives:

$$abla_{ heta} J(heta) = \sum_{s=0}^{n}
abla_{a} Q_{\phi}(s, a) \cdot
abla_{ heta} \pi_{ heta}(s)$$

Interpretation:

- ▶ Train $Q_{\phi}(s,a)$ from experience
- Update policy in the direction of actions with higher value
- ► 'Let the policy follow the value network'

Motivation for DDPG

Goal: Extend value-based deep RL (like DQN) to **continuous action spaces**. Challenge in continuous spaces:

$$a^*(s) = \arg\max_{a} Q^*(s, a)$$
 is hard!

DDPG solution:

- ▶ Use the gradient of Q(s, a) w.r.t. the action to approximate max_a Q(s, a)
- ► Deterministic policy actor: $\pi_{\theta}(s)$ ► Critic estimates $Q_{\phi}(s, a)$

Based on:

- ► Deterministic policy gradient²
- ► NFQCA³

Conference on Machine Learning. 2014, pp. 387–395.

³Roland Hafner and Martin Riedmiller. 'Reinforcement learning in feedback control'. In:

Machine Learning 84.1-2 (2011), pp. 137-169.

²David Silver et al. 'Deterministic policy gradient algorithms'. In: *International*

DDPG Algorithm Pseudocode

Randomly initialize critic network $Q_{\phi}(s,a)$ and actor $\pi_{\theta}(s)$ with weights ϕ and θ . Initialize target network Q' and π' with weights $\phi' \leftarrow \phi$, $\theta' \leftarrow \theta$ Initialize replay buffer R for episode $=1\ldots M$ do

Initialize a random process $\mathcal N$ for action exploration

Receive initial observation state s_1

for $t = 1 \dots T$ do

 $a_t = \pi_{\theta}(s_t) + \mathcal{N}_t$ according to the current policy and exploration noise Execute action a_t and observe reward r_t and observe new state s_{t+1}

Store transition (s_t, a_t, r_t, s_{t+1}) in R

Sample a random minibatch of N transitions (s_i, a_i, r_i, s_{i+1}) from R

Set $y_i = r_i + \gamma Q_{\phi'}(s_{i+1}, \pi_{\theta'}(s_{i+1}))$

Update critic by minimizing the loss: $L = \frac{1}{N} \sum_{i} (y_i - Q_{\phi}(s_i, a_i))^2$

Update the actor policy using the sampled policy gradient:

$$abla_{ heta} J pprox rac{1}{N} \sum_{i}
abla_{a} Q_{\phi}(s,a) |_{s=s_{i},a=\mu(s_{i})}
abla_{ heta} \pi_{ heta}(s) |_{s_{i}}$$

DDPG Algorithm Pseudocode

Update the target networks:

$$\phi' \leftarrow \tau \phi + (1 - \tau)\phi'$$
$$\theta' \leftarrow \tau \theta + (1 - \tau)\theta'$$

end for end for

DDPG Algorithm Overview

DDPG characteristics:

- Actor-critic, model-free
- Continuous actions
- ► Off-policy learning
- ► Replay buffer (like DQN)
- ► Target networks for stability (like DQN)

DDPG showed strong performance on physics control tasks:

- ► CartPole, Gripper, Walker, Car driving
- ► Learns from pixels in some settings

Key idea: Use deterministic policy + learned $\it Q$ gradient to update actor.

DDPG Resources

Implementations:

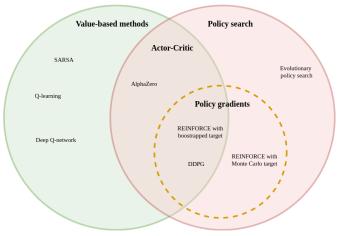
- Spinning Up (spinningup.openai.com)
- Stable-Baselines (stable-baselines.readthedocs.io)
- ► Original paper⁴

DDPG Summary:

- Bridges DQN ideas into continuous control
- Strong baseline in robotics and physics simulators
- ► Foundation for later methods (TD3, SAC)

⁴Timothy P Lillicrap et al. 'Continuous control with deep reinforcement learning'. In: International Conference on Learning Representations. 2016.

Value vs Policy vs Actor-Critic



Value-based, policy-based, and actor-critic methods⁵.

⁵Thomas Moerland. *Continuous Markov Decision Process and Policy Search*. Lecture Notes for the Course Reinforcement Learning. Leiden University, 2021.