# Actor-Critic

Hung-yi Lee

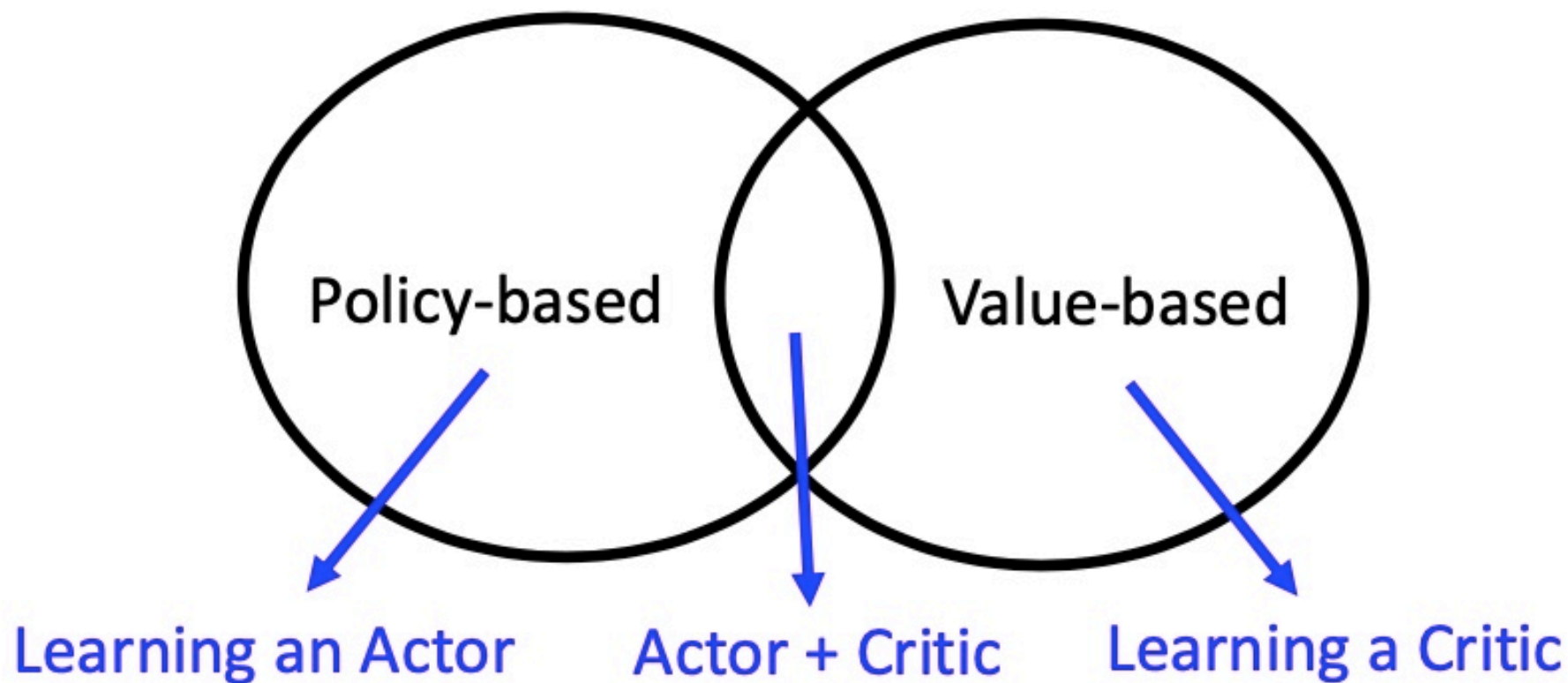Policy-based     Value-based

Learning an Actor     Actor + Critic     Learning a Critic
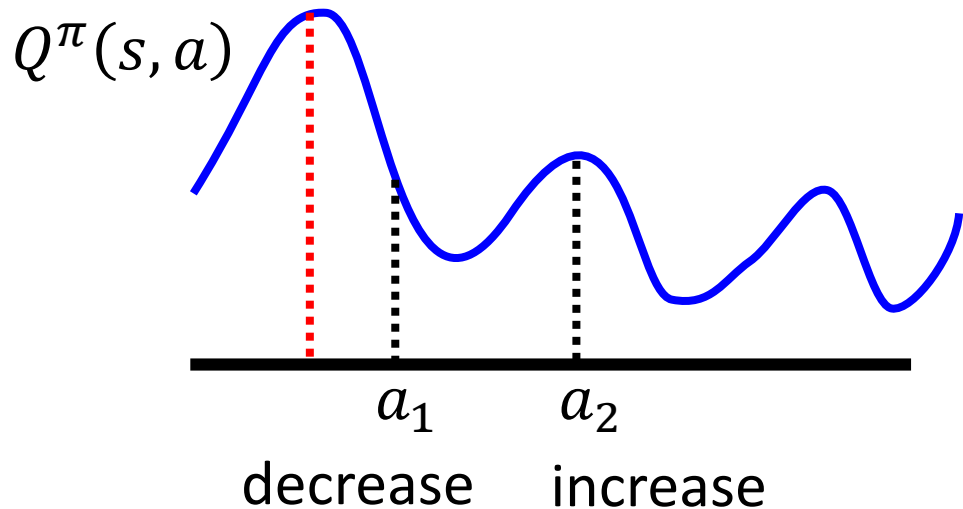
# Pathwise Derivative Policy Gradient

David Silver, Guy Lever, Nicolas Heess, Thomas Degris, Daan Wierstra, Martin Riedmiller, "Deterministic Policy Gradient Algorithms", ICML, 2014

Timothy P. Lillicrap, Jonathan J. Hunt, Alexander Pritzel, Nicolas Heess,
Tom Erez, Yuval Tassa, David Silver, Daan Wierstra, "CONTINUOUS CONTROL WITH DEEP REINFORCEMENT LEARNING", ICLR, 2016
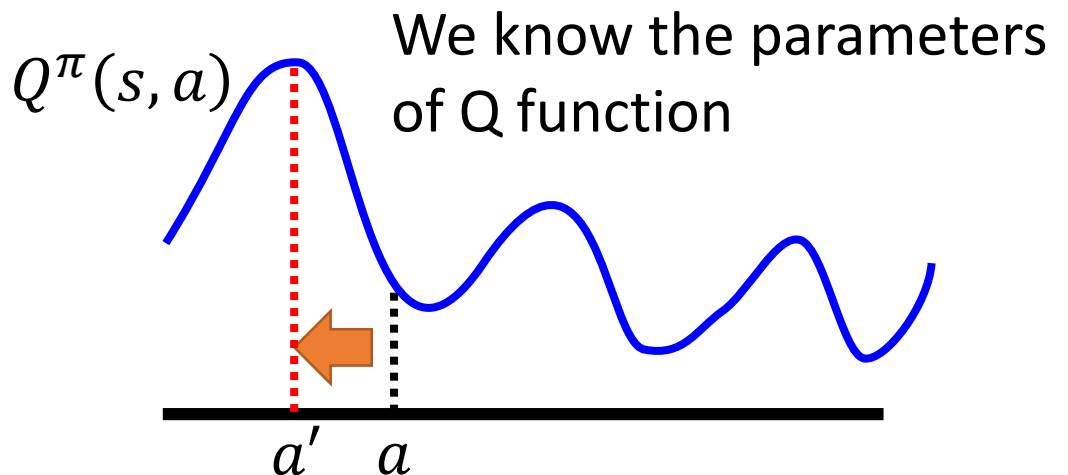
# Another Way to use Critic

**_Original Actor-critic_**

$Q^\pi(s,a)$

$a_1$    $a_2$

decrease    increase

**_Pathwise derivative policy gradient_**

From Q function we know that taking a' at state s is better than a

We know the parameters of Q function

$Q^\pi(s,a)$

$a'$    $a$

# Actor



# Critic

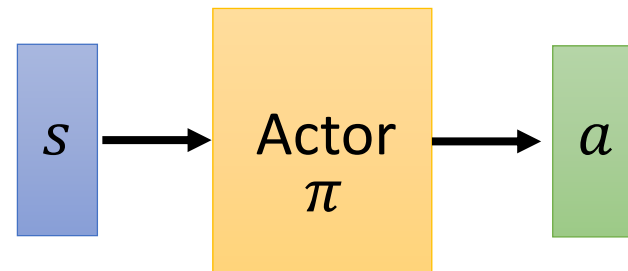

Pathwise derivative policy gradient

Original Actor-critic

---

Action $a$ is a *continuous vector*

$$a = arg \max_a Q(s, a)$$

$$s \rightarrow \boxed{\text{Actor } \pi} \rightarrow a$$

Actor as the solver of this optimization problem

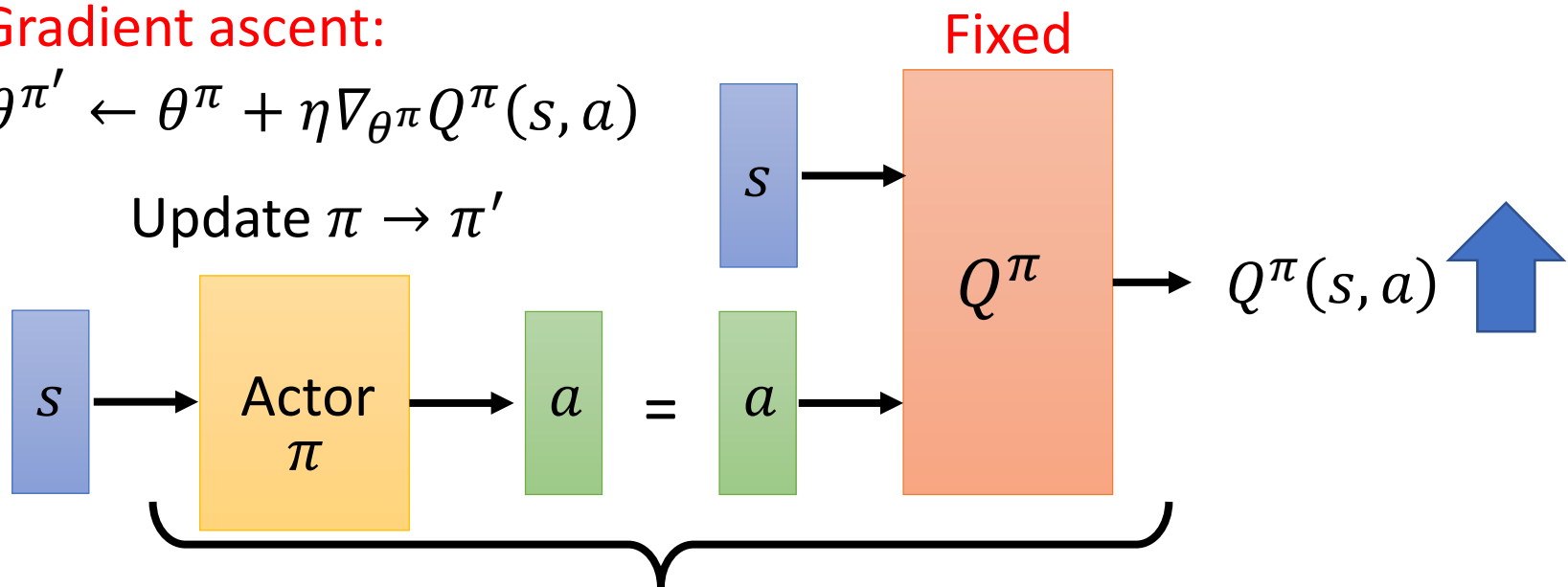# Pathwise Derivative Policy Gradient
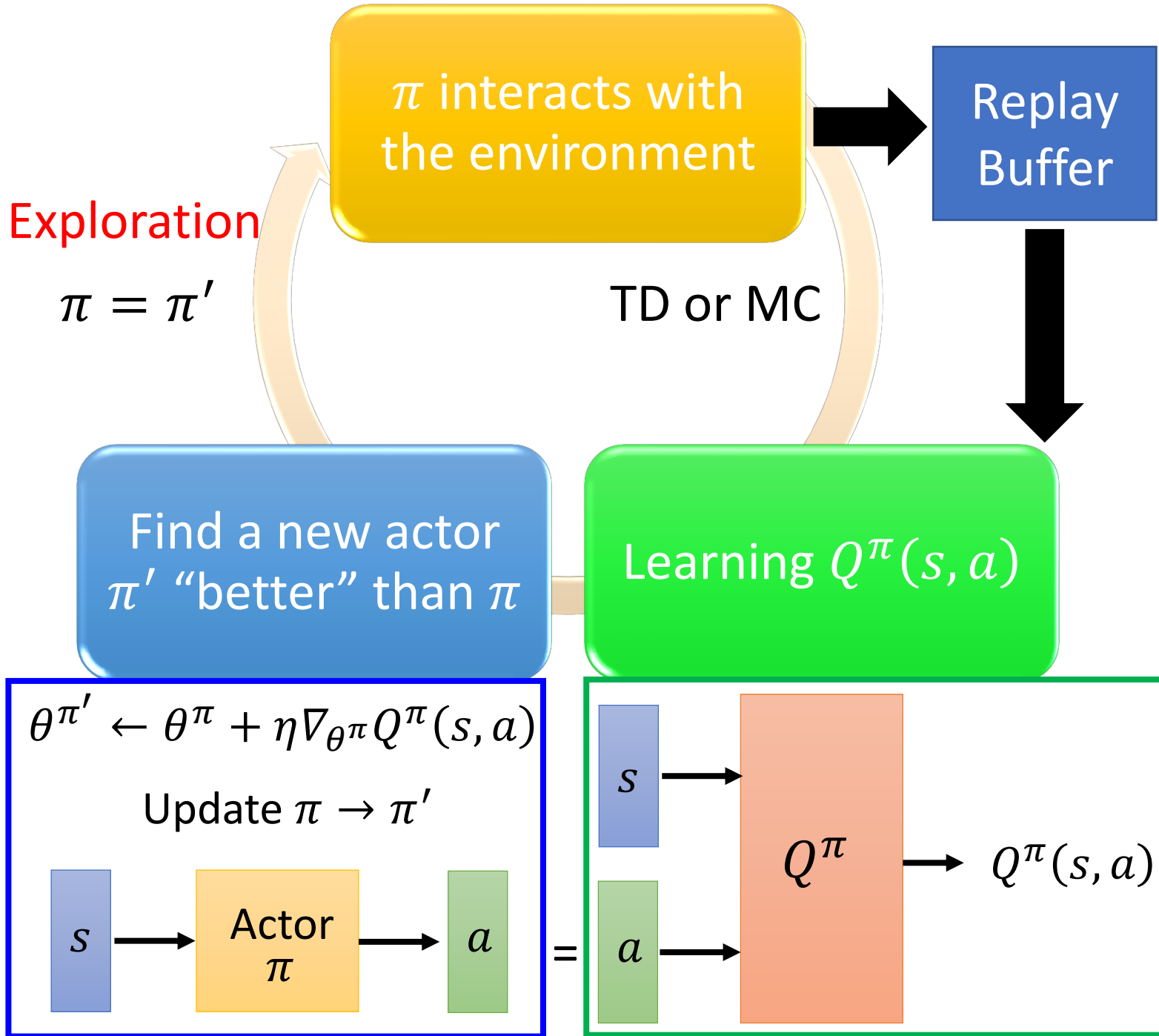
$$\pi'(s) = arg \max_a Q^\pi(s,a)$$ ⬅ a is the output of an actor

Gradient ascent:
$$\theta^{\pi'} \leftarrow \theta^\pi + \eta \nabla_{\theta^\pi} Q^\pi(s,a)$$

Update $\pi \to \pi'$

Fixed



$s$ → Actor $\pi$ → $a$ = $a$ → $Q^\pi$ → $Q^\pi(s,a)$ ⬆

This is a large network

## *Q-Learning Algorithm*

- Initialize Q-function $Q$, target Q-function $\hat{Q} = Q$

- In each episode
  - For each time step t
    - Given state $s_t$, take action $a_t$ based on Q (exploration)
    - Obtain reward $r_t$, and reach new state $s_{t+1}$
    - Store $(s_t, a_t, r_t, s_{t+1})$ into buffer
    - Sample $(s_i, a_i, r_i, s_{i+1})$ from buffer (usually a batch)
    - Target $y = r_i + \max_{a} \hat{Q}(s_{i+1}, a)$
    - Update the parameters of $Q$ to make $Q(s_i, a_i)$ close to $y$ (regression)

    - Every C steps reset $\hat{Q} = Q$

# *Q-Learning Algorithm* ➡ *Pathwise Derivative Policy Gradient*

- Initialize Q-function $Q$, target Q-function $\hat{Q} = Q$, actor $\pi$, target actor $\hat{\pi} = \pi$

- In each episode
  - For each time step t
  - **1** • Given state $s_t$, take action $a_t$ based on ~~$Q$~~ $\pi$ (exploration)
    - Obtain reward $r_t$, and reach new state $s_{t+1}$
    - Store $(s_t, a_t, r_t, s_{t+1})$ into buffer
    - Sample $(s_i, a_i, r_i, s_{i+1})$ from buffer (usually a batch)
  - **2** • Target $y = r_i + \underset{a}{\max} \hat{Q}(s_{i+1}, a) \hat{Q}\big(s_{i+1}, \hat{\pi}(s_{i+1})\big)$
    - Update the parameters of $Q$ to make $Q(s_i, a_i)$ close to $y$ (regression)
  - **3** • Update the parameters of $\pi$ to maximize $Q\big(s_i, \pi(s_i)\big)$
    - Every C steps reset $\hat{Q} = Q$
  - **4** • Every C steps reset $\hat{\pi} = \pi$

# Connection with GAN

| Method | GANs | AC |
|---|---|---|
| Freezing learning | yes | yes |
| Label smoothing | yes | no |
| Historical averaging | yes | no |
| Minibatch discrimination | yes | no |
| Batch normalization | yes | yes |
| Target networks | n/a | yes |
| Replay buffers | no | yes |
| Entropy regularization | no | yes |
| Compatibility | no | yes |

David Pfau, Oriol Vinyals, "Connecting Generative Adversarial Networks and Actor-Critic Methods", arXiv preprint, 2016