# Introduction of Reinforcement Learning

# Machine Learning ≈Looking for a Func



Observation

**Function input**

Actor/Policy

Action =
π( Observation )

Action

**Function output**

**Used to pick the best function**

Reward

Environment
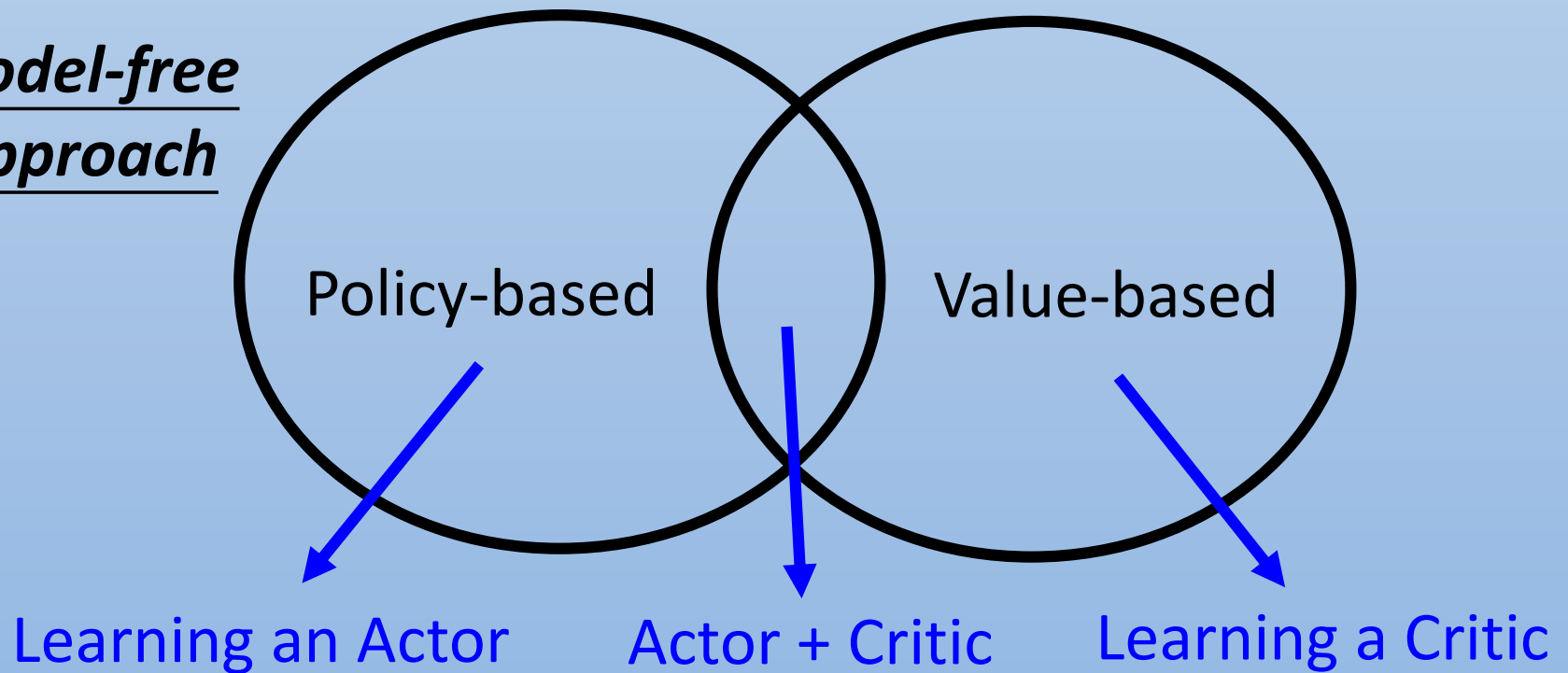
# Properties of Reinforcement Learning

- **Reward delay**
  - In space invader, only "fire" obtains reward
    - Although the moving before "fire" is important
  - In Go playing, it may be better to sacrifice immediate reward to gain more long-term reward
- Agent's actions **affect the subsequent data it receives**
  - E.g. Exploration

# Outline

Alpha Go: policy-based + value-based + model-based
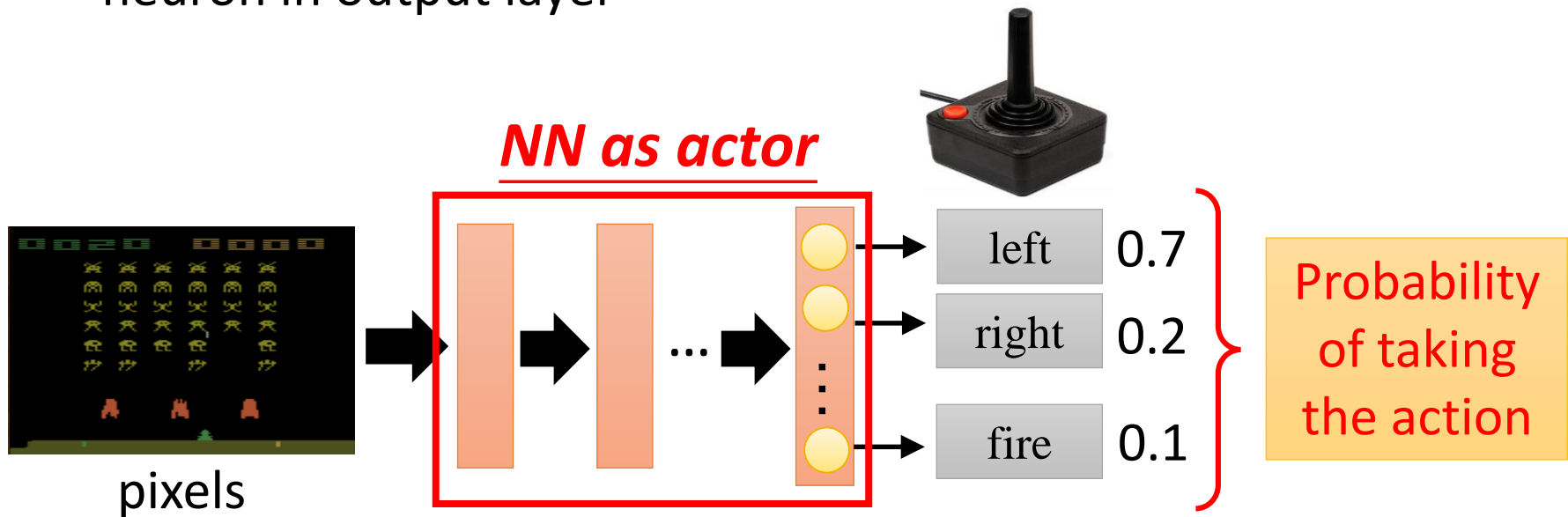
***Model-free Approach***

Policy-based

Value-based

Learning an Actor

Actor + Critic

Learning a Critic

***Model-based Approach***

# Policy-based Approach

Learning an Actor

# Neural network as Actor

- Input of neural network: the observation of machine represented as a vector or a matrix

- Output neural network : each action corresponds to a neuron in output layer

**NN as actor**



left 0.7

right 0.2

fire 0.1

pixels

Probability of taking the action

What is the benefit of using network instead of lookup table?

generalization

# Goodness of Actor

- Given an actor $\pi_\theta(s)$ with network parameter $\theta$
- Use the actor $\pi_\theta(s)$ to play the video game

We define $\bar{R}_\theta$ as the _expected value_ of $R_\theta$
$\bar{R}_\theta$ evaluates the goodness of an actor $\pi_\theta(s)$

- An episode is considered as a trajectory $\tau$
  - $\tau = \{s_1, a_1, r_1, s_2, a_2, r_2, \cdots, s_T, a_T, r_T\}$
  - $R(\tau) = \sum_{t=1}^{T} r_t$
  - If you use an actor to play the game, each $\tau$ has a probability to be sampled
    - The probability depends on actor parameter $\theta$: $P(\tau|\theta)$

Total reward: $R_\theta = \sum_{t=1}^{T} r_t,$

Even with the same actor

$R_\theta$ is different each time

$$\bar{R}_\theta = \sum_\tau R(\tau)P(\tau|\theta) \approx \frac{1}{N}\sum_{n=1}^{N} R(\tau^n)$$

Sum over all possible trajectory

Use $\pi_\theta$ to play the game N times,
obtain $\{\tau^1, \tau^2, \cdots, \tau^N\}$

Sampling $\tau$ from $P(\tau|\theta)$ N times

# Gradient Ascent

- Problem statement

$$\theta^* = arg \max_\theta \bar{R}_\theta$$

- Gradient ascent
  - Start with $\theta^0$
  - $\theta^1 \leftarrow \theta^0 + \eta \nabla \bar{R}_{\theta^0}$
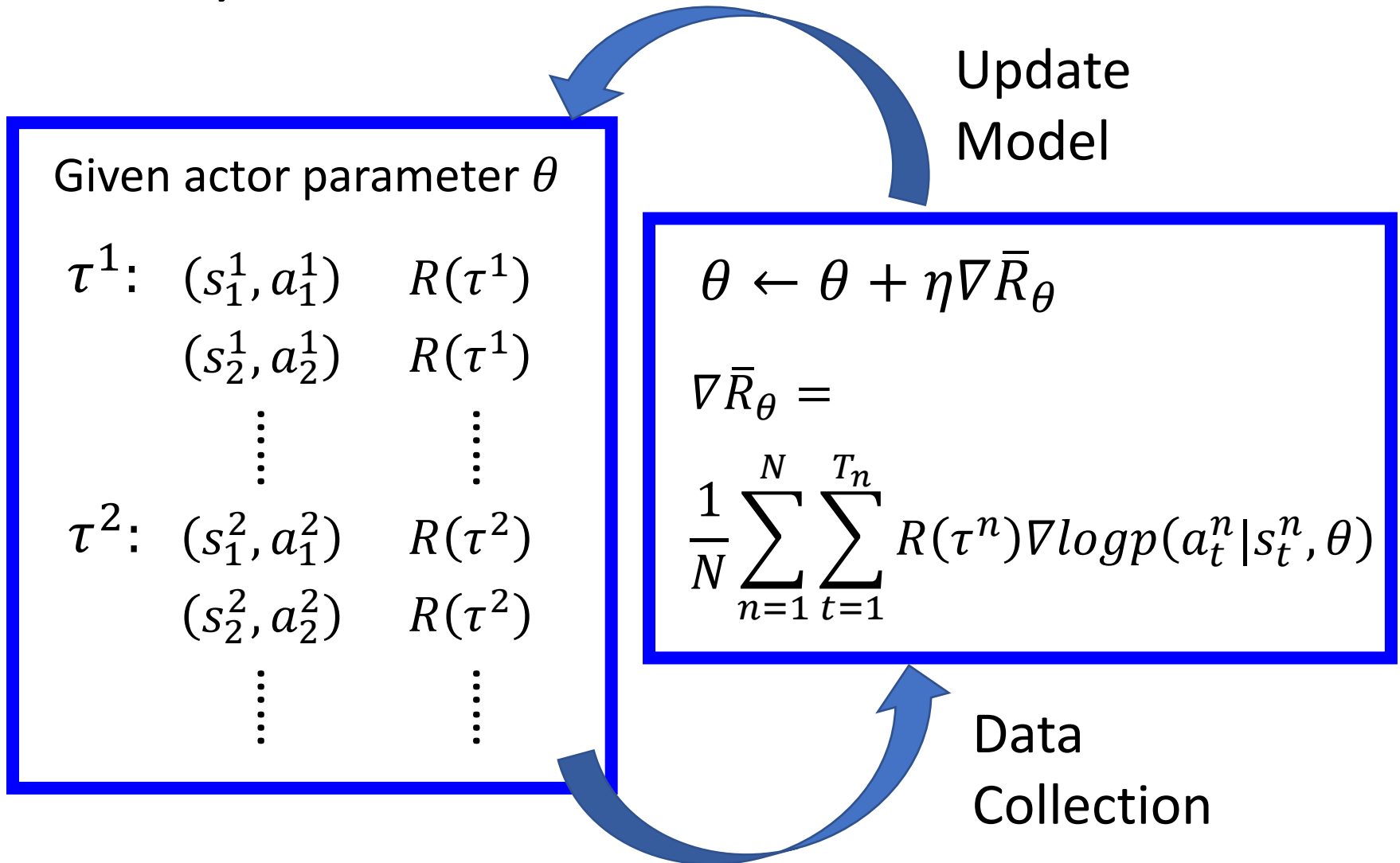  - $\theta^2 \leftarrow \theta^1 + \eta \nabla \bar{R}_{\theta^1}$
  - ……

$$\theta = \{w_1, w_2, \cdots, b_1, \cdots\}$$

$$\nabla \bar{R}_\theta = \begin{bmatrix} \partial \bar{R}_\theta / \partial w_1 \\ \partial \bar{R}_\theta / \partial w_2 \\ \vdots \\ \partial \bar{R}_\theta / \partial b_1 \\ \vdots \end{bmatrix}$$
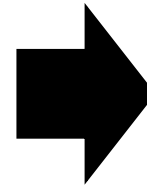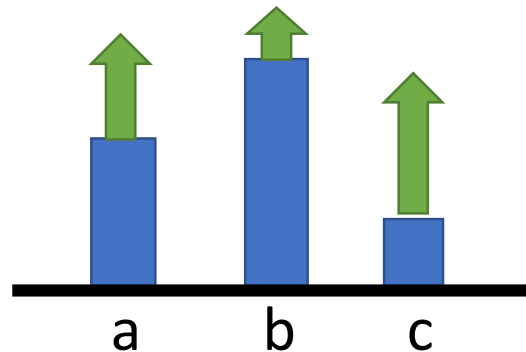
# Policy Gradient

Given actor parameter $\theta$

$\tau^1$:    $(s_1^1, a_1^1)$    $R(\tau^1)$

       $(s_2^1, a_2^1)$    $R(\tau^1)$

         $\vdots$        $\vdots$

$\tau^2$:    $(s_1^2, a_1^2)$    $R(\tau^2)$

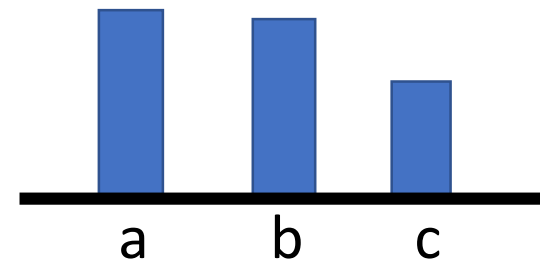       $(s_2^2, a_2^2)$    $R(\tau^2)$

         $\vdots$        $\vdots$

Update Model

$$\theta \leftarrow \theta + \eta \nabla \bar{R}_\theta$$

$$\nabla \bar{R}_\theta =$$

$$\frac{1}{N} \sum_{n=1}^{N} \sum_{t=1}^{T_n} R(\tau^n) \nabla log p(a_t^n | s_t^n, \theta)$$

Data Collection

# Add a Baseline

It is possible that $R(\tau^n)$ is always positive.

$$\theta^{new} \leftarrow \theta^{old} + \eta \nabla \bar{R}_{\theta^{old}}$$

$$\nabla \bar{R}_\theta \approx \frac{1}{N} \sum_{n=1}^{N} \sum_{t=1}^{T_n} (R(\tau^n) - b) \nabla log p(a_t^n | s_t^n, \theta)$$

Ideal case
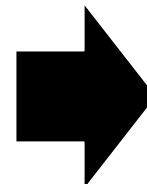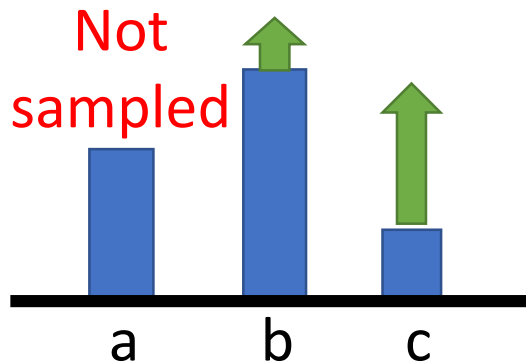
It is probability …

a   b   c

a   b   c

Sampling
……

Not sampled

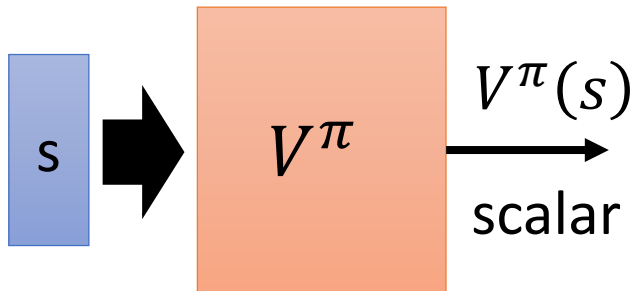The probability of the actions not sampled will decrease.

a   b   c

a   b   c

# Value-based Approach

## Learning a Critic

# Critic

- State value function $V^\pi(s)$
  - When using actor $\pi$, the *cumulated* reward expects to be obtained after seeing observation (state) s
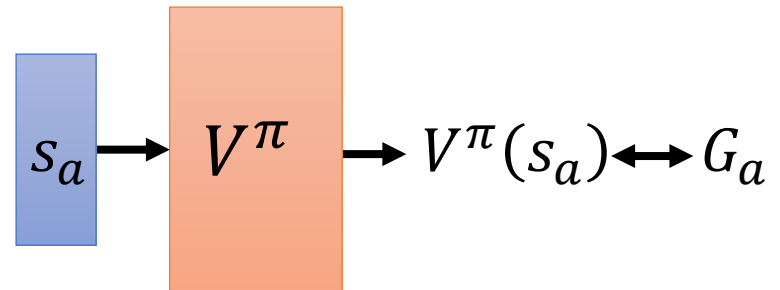


$V^\pi(s)$ is large

$V^\pi(s)$ is smaller

# How to estimate $V^\pi(s)$

- Monte-Carlo based approach
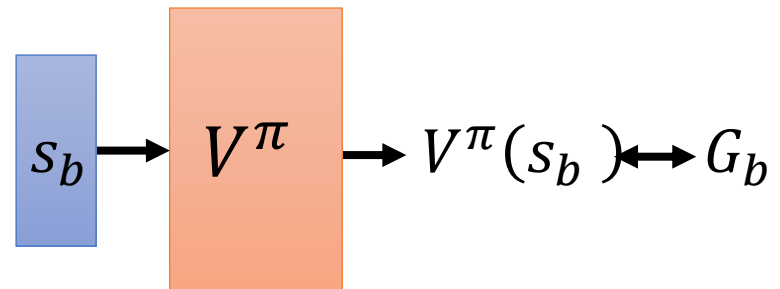  - The critic watches $\pi$ playing the game

After seeing $s_a$,

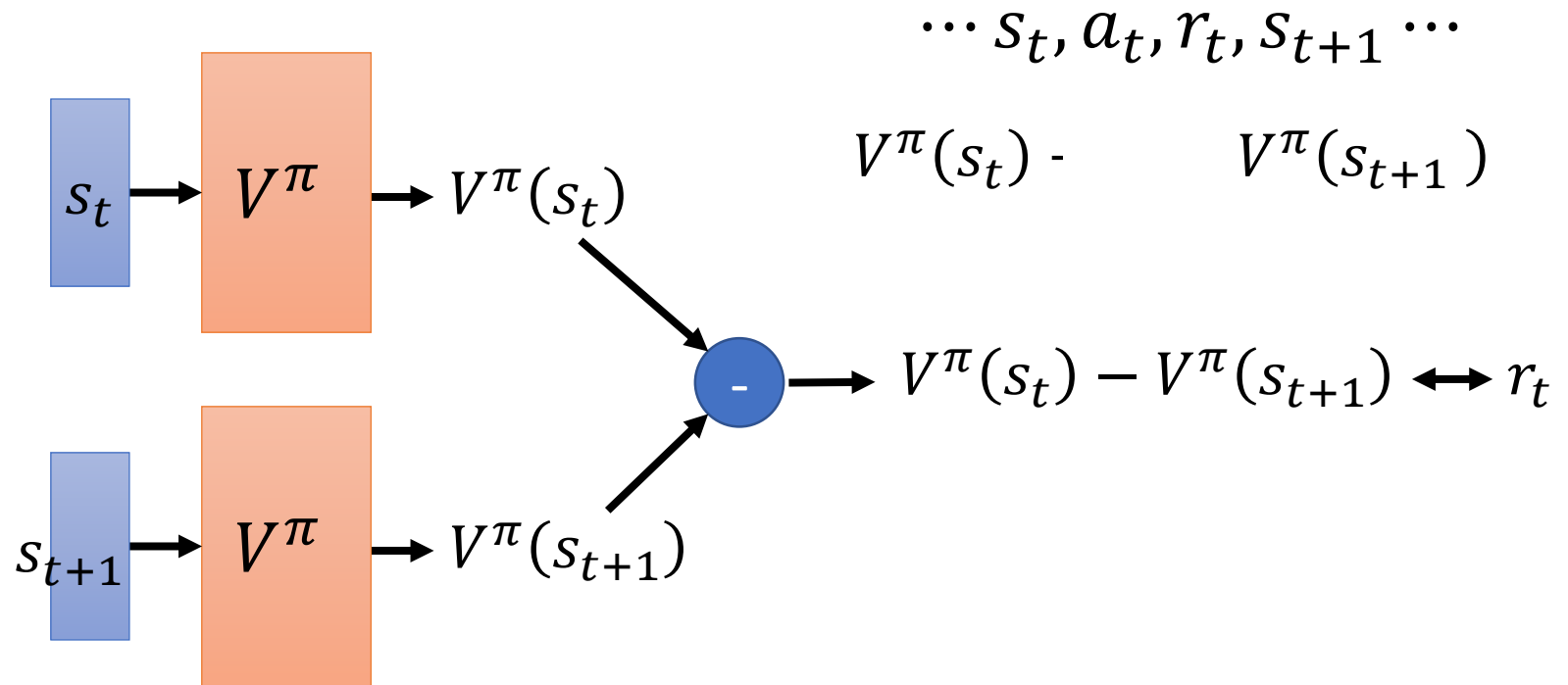Until the end of the episode, the cumulated reward is $G_a$

$$s_a \rightarrow V^\pi \rightarrow V^\pi(s_a) \longleftrightarrow G_a$$

After seeing $s_b$,

Until the end of the episode, the cumulated reward is $G_b$

$$s_b \rightarrow V^\pi \rightarrow V^\pi(s_b) \longleftrightarrow G_b$$

# How to estimate $V^\pi(s)$

- Temporal-difference approach



$$\cdots s_t, a_t, r_t, s_{t+1} \cdots$$

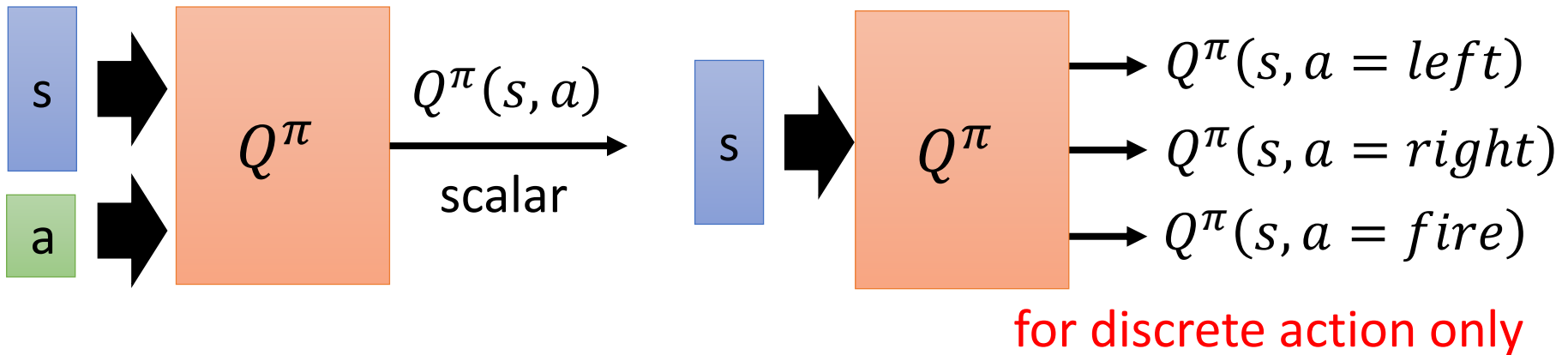$$V^\pi(s_t) - \qquad V^\pi(s_{t+1})$$

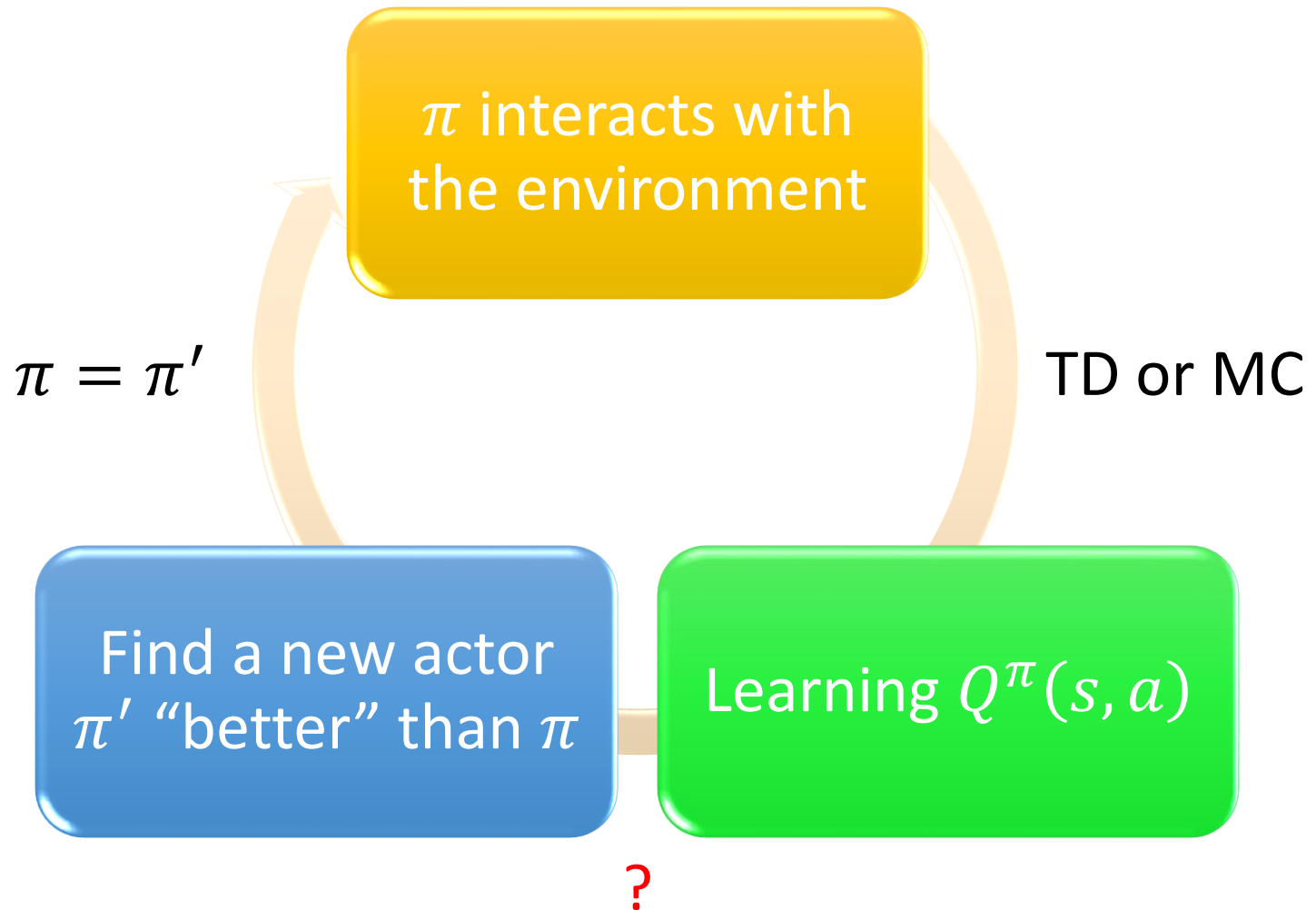$$V^\pi(s_t) - V^\pi(s_{t+1}) \longleftrightarrow r_t$$

Some applications have very long episodes, so that delaying all learning until an episode's end is too slow.

# Another Critic

- State-action value function $Q^\pi(s, a)$
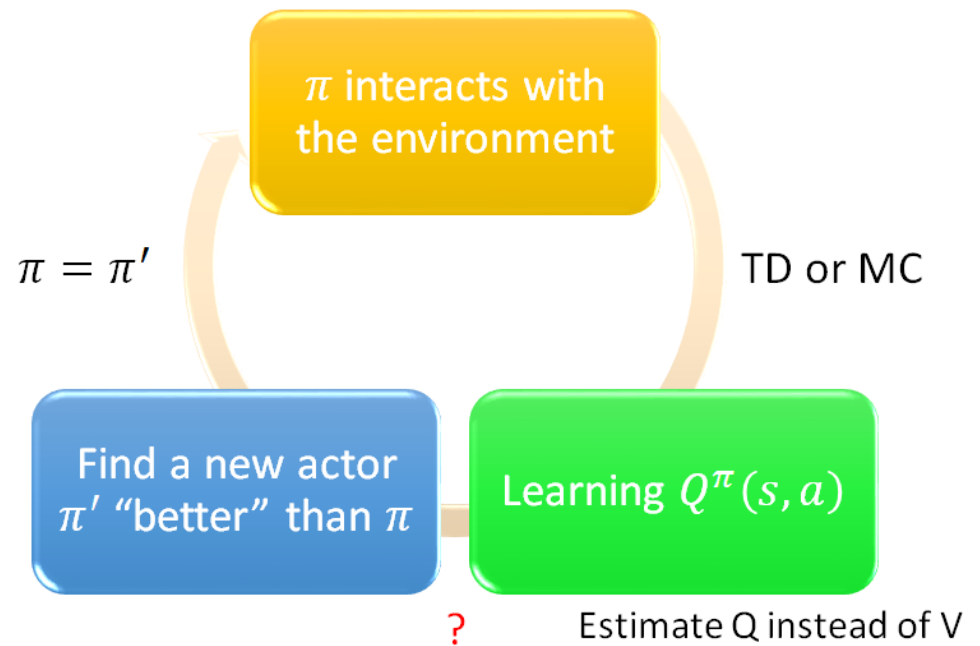  - When using actor $\pi$, the *cumulated* reward expects to be obtained after seeing observation s and taking a



for discrete action only

# Q-Learning

# Q-Learning



- Given $Q^{\pi}(s, a)$, find a new actor $\pi'$ "better" than $\pi$
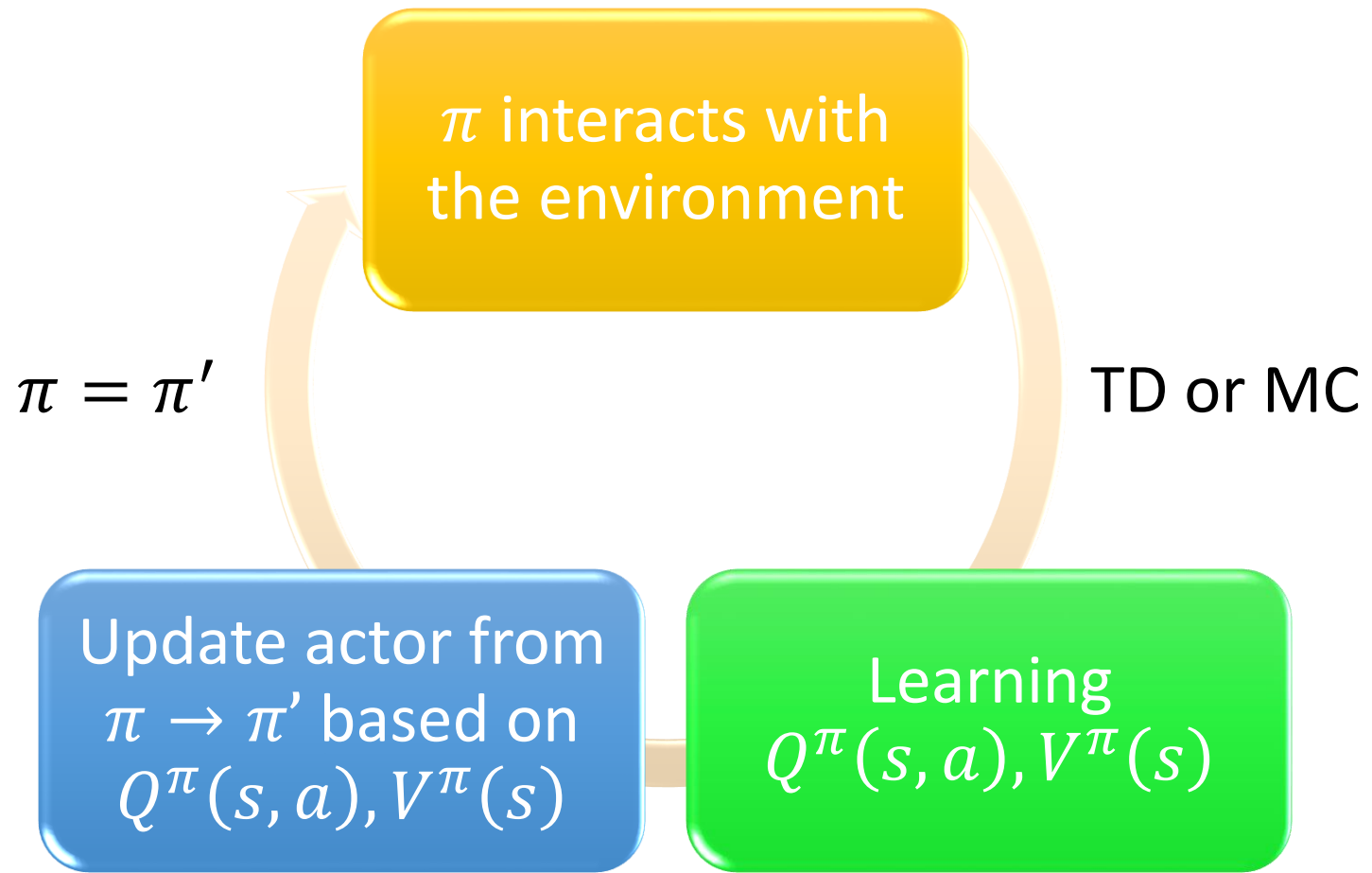  - "Better": $V^{\pi'}(s) \geq V^{\pi}(s)$, for all state s

$$\pi'(s) = arg \max_{a} Q^{\pi}(s, a)$$

➢ $\pi'$ does not have extra parameters. It depends on Q

➢ Not suitable for continuous action a

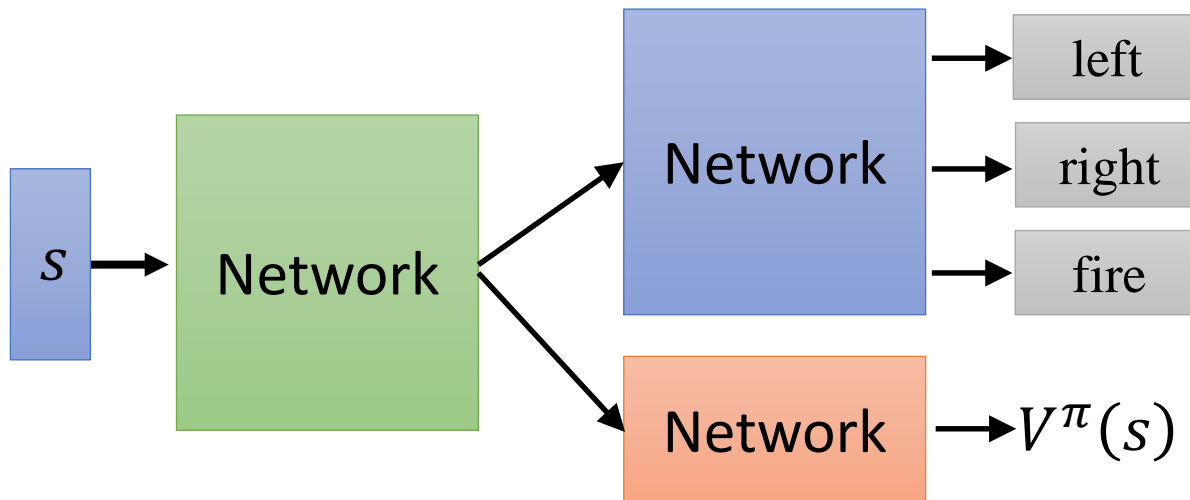# Deep Reinforcement Learning

Actor-Critic

# Actor-Critic

# Actor-Critic

- Tips
  - The parameters of actor $\pi(s)$ and critic $V^{\pi}(s)$ can be shared

# *Asynchronous*

Source of image:
https://medium.com/emergent-future/simple-reinforcement-learning-with-tensorflow-part-8-asynchronous-actor-critic-agents-a3c-c88f72a5e9f2#.68x6na7o9

1. Copy global parameters

2. Sampling some data

3. Compute gradients

4. Update global models



$$\cancel{\theta^1} + \eta \Delta \theta$$

$$\theta^2$$

(other workers also update models)

$$\Delta \theta$$

$$\theta^1$$

$$\theta^1$$

$$\Delta \theta$$