

A Brief Introduction to Reinforcement Learning

Minlie Huang (黄民烈)

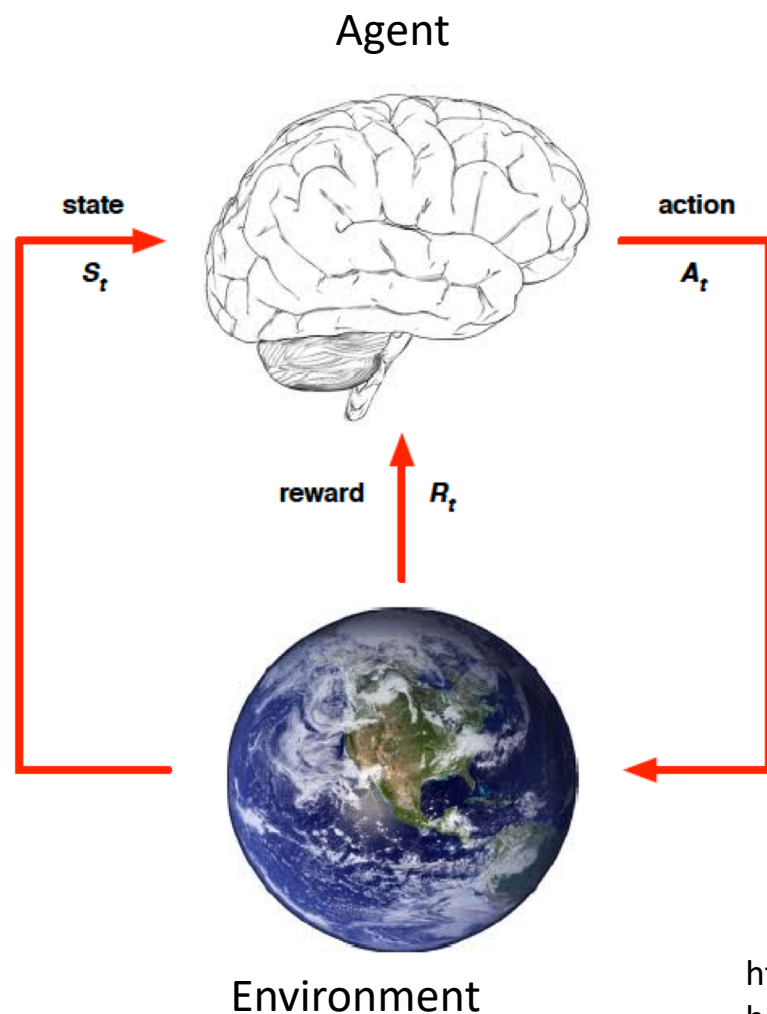
Dept. of Computer Science,
Tsinghua University

aihuang@tsinghua.edu.cn

<http://coai.cs.tsinghua.edu.cn/hml>



Reinforcement Learning

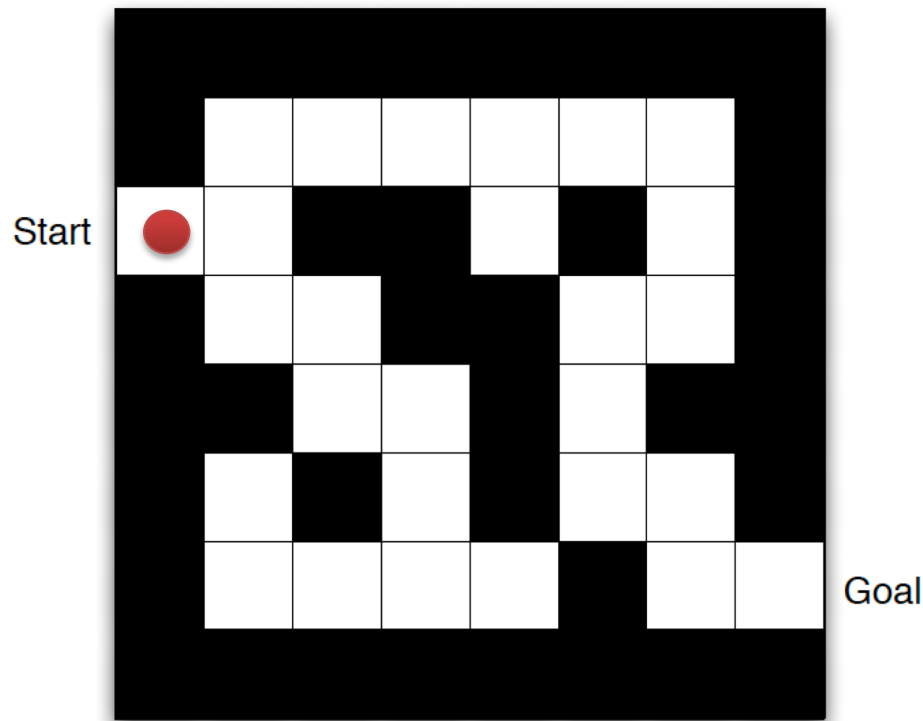


At each step t :

- The agent receives a **state** S_t from environment
- The agent executes **action** A_t based on the received state
- The agent receives scalar **reward** R_t from the environment
- The environment transform into a new state S_{t+1}



Maze Example



States: Agent's location

Actions: N, E, S, W

Rewards:

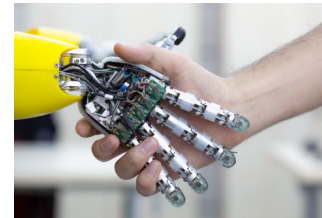
- 100 if reaching the goal
- -100 if reaching the dead end
- -1 per time-step



Deep Reinforcement Learning



Deep learning to represent **states**, **actions**,
or **policy functions**



Robotics, control



Self-driving



Language interaction



System operating



Reinforcement Learning

◉ Markov Decision Process (MDP)

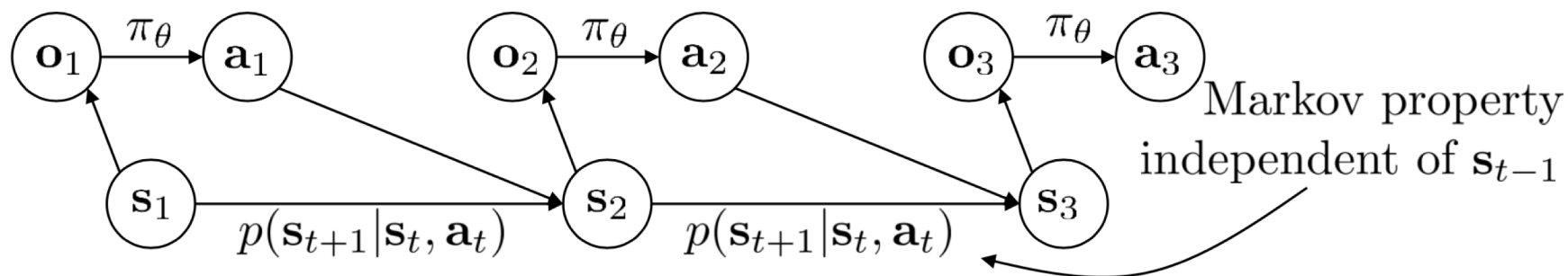
\mathbf{s}_t – state

\mathbf{o}_t – observation

\mathbf{a}_t – action

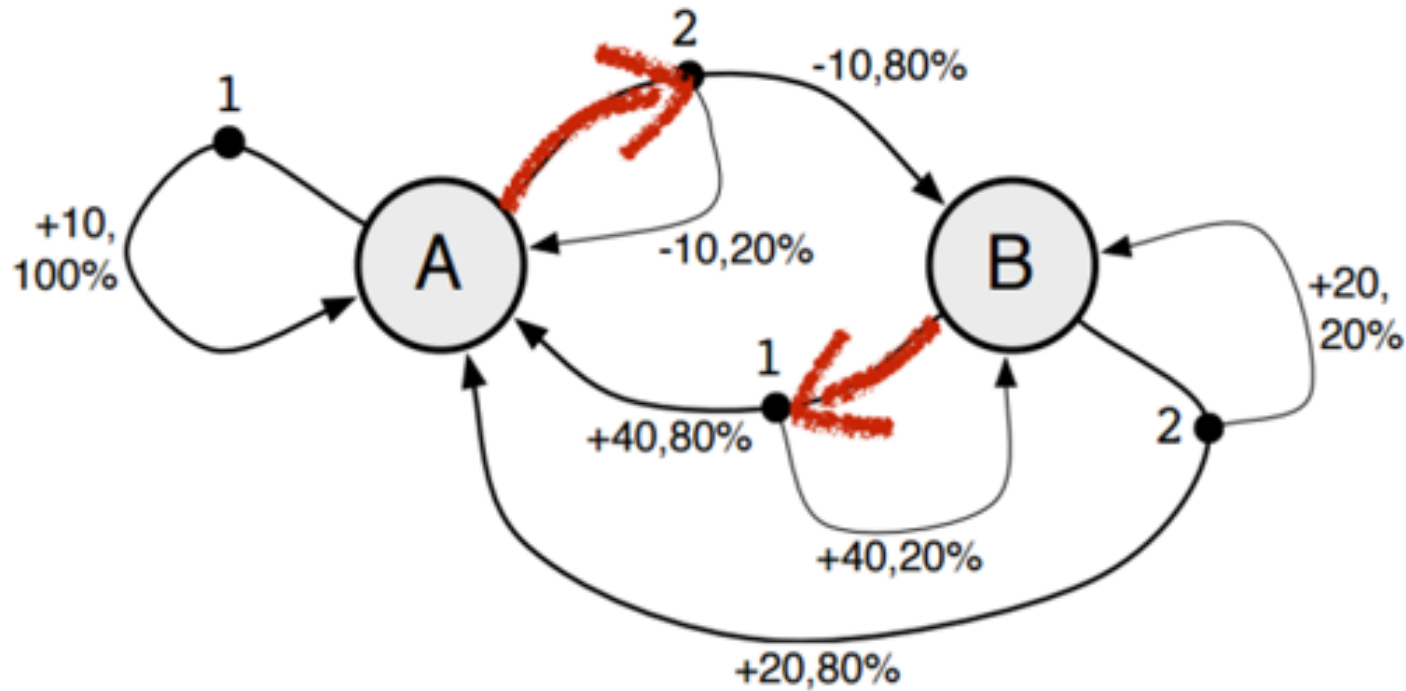
$\pi_\theta(\mathbf{a}_t|\mathbf{o}_t)$ – policy

$\pi_\theta(\mathbf{a}_t|\mathbf{s}_t)$ – policy (fully observed)



A MDP Game

If the agent is at state A, which action would you prefer?



Reinforcement Learning

- Policy: Agent's behavior strategy
 - How to choose an action given the current state in order to maximize the total reward
- Deterministic policy: $a = \pi(o)$
- Stochastic policy: $a \sim \pi(o) = p(a|o)$
- Fully observable: $o = s$
- Partially observable: $o \neq s$



Value-Based RL

Value Based

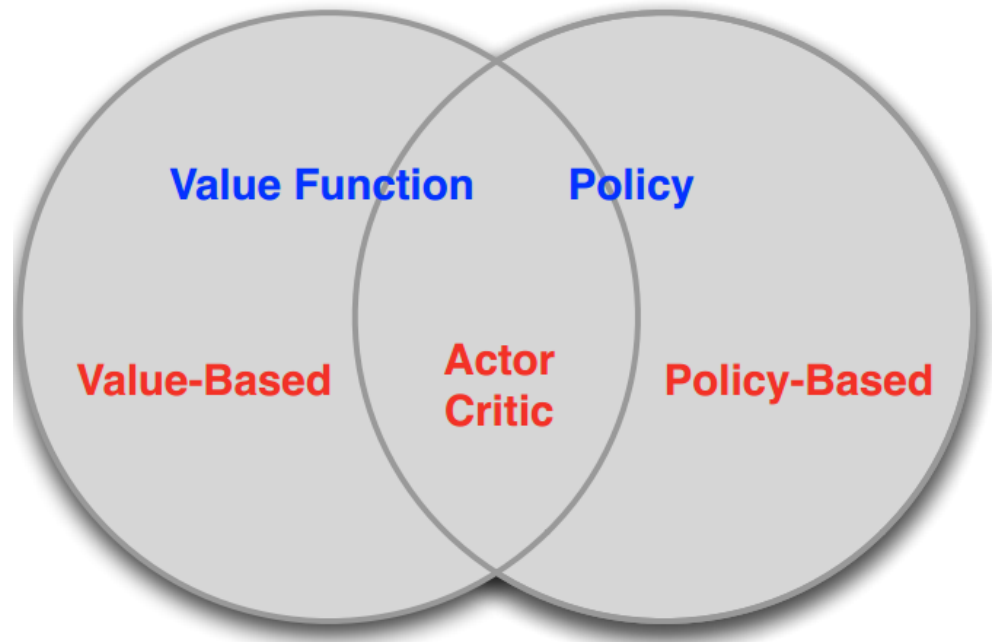
- ◆ Learnt Value Function
- ◆ Implicit Policy

Policy Based

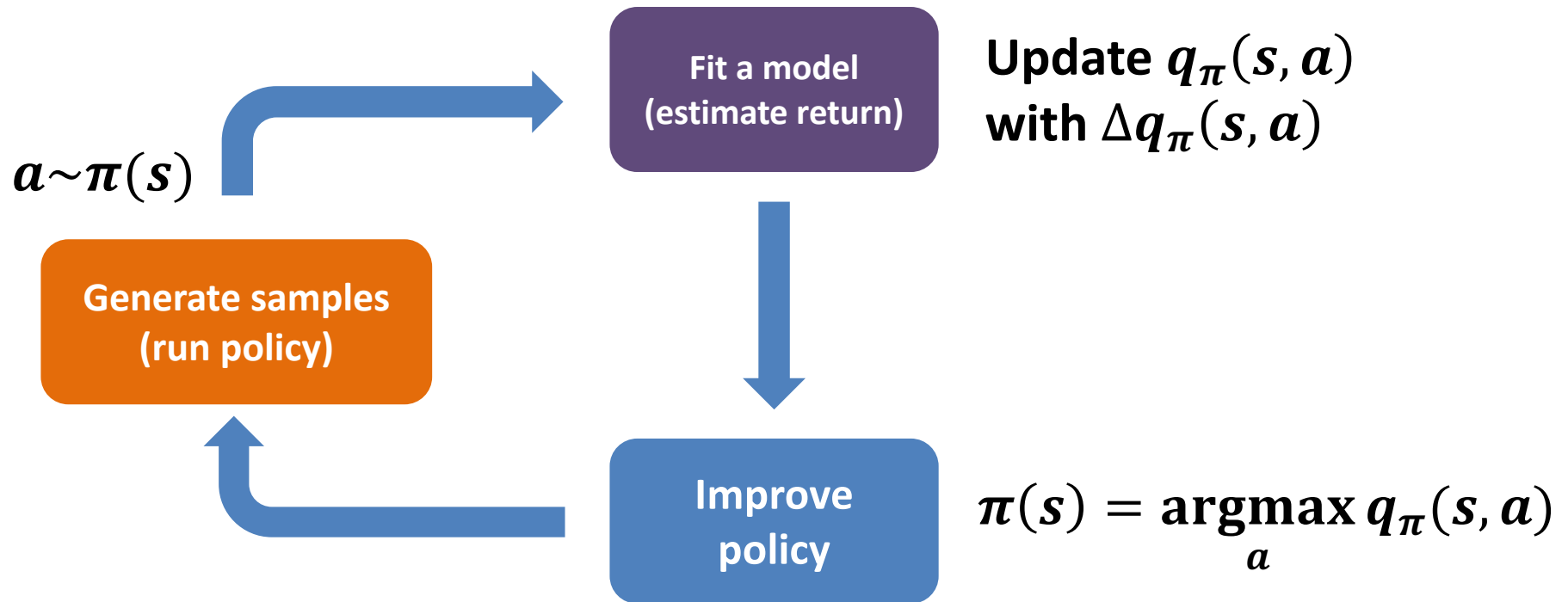
- ◆ No Value Function
- ◆ Learnt Policy

Actor-Critic

- ◆ Learnt Value Function
- ◆ Learnt Policy



Q-Learning: Procedure



Value Function

- ◎ A **value function** is a prediction of **future reward**
 - ◆ How much reward will I get from **action a** given **state s** under policy **π** ?
- ◎ **Q-value function** (action-value function) gives expected total reward
 - ◆ Each state-action pair **(s, a)** has an entry **$q(s, a)$**



Value Function

◎ **Q-value function** gives expected total reward

- ◆ Sample an trajectory with policy π
- ◆ ... $S_t, A_t, R_{t+1}, S_{t+1}, A_{t+1}, R_{t+2}, S_{t+2}, \dots$
- ◆ Discount factor γ : between $[0, 1]$, **how far ahead in time** the algorithm looks

$$q_{\pi}(s, a) = \mathbb{E} \left[\underbrace{R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots}_{\text{Delayed reward is taken into consideration}} \mid S_t = s, A_t = a, A_{t+1:\infty} \sim \pi \right]$$

Delayed reward is taken into consideration

◎ **Future rewards!**

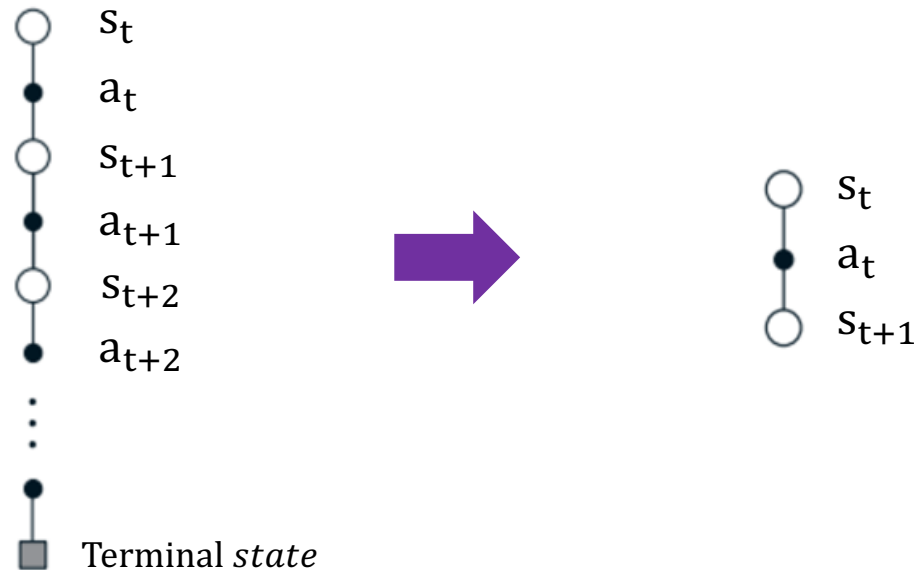


Value Function

- ◎ Bellman equation (with a particular policy π)

$$q_{\pi}(s, a) = \mathbb{E} \left[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \cdots \mid S_t = s, A_t = a, A_{t+1:\infty} \sim \pi \right]$$

$$q_{\pi}(s, a) = \mathbb{E} \left[R_{t+1} + \gamma q_{\pi}(S_{t+1}, A_{t+1}) \mid S_t = s, A_t = a, A_{t+1} \sim \pi \right]$$



Optimal Value Function

- ⊙ An optimal value function is the maximum achievable value

$$q^*(s, a) = \max_{\pi} [q_{\pi}(s, a)] = q_{\pi^*}(s, a)$$

- ⊙ Act optimally:

$$\pi^*(s) = \operatorname{argmax}_a [q^*(s, a)]$$

- ⊙ Optimal value maximizes over all decisions

$$q^*(s, a) = R_{t+1} + \gamma \max_{a_{t+1}} R_{t+2} + \gamma^2 \max_{a_{t+2}} R_{t+3}$$

$$= R_{t+1} + \gamma \max_{a_{t+1}} q^*(s_{t+1}, a_{t+1})$$



Q-Learning: Optimization

- Optimal q-values obey Bellman function

$$q^*(s, a) = E_{s'}[R + \gamma \max_{a'} q^*(s', a') | s, a]$$

- Minimizing mean-squared error of **value function** between *approximate value* and *true value (target)*

$$\Delta q(s, a) = \underbrace{\alpha(R + \gamma \max_{a'} q(s', a'))}_{\text{Biased true value}} - \underbrace{q(s, a)}_{\text{Approximate value}}$$

- Policy $\pi(s)$: choose the action that has the maximum **profit**



Q-Learning: Trial and Error

⊙ Exploration & Exploitation : ϵ -greedy

◆ Require sufficient exploration to enable learning

$$\text{◆ } a_t = \begin{cases} \arg \max_a q(s_t, a) & \text{with probability } 1 - \epsilon \\ \text{random action} & \text{with probability } \epsilon \end{cases}$$



Q-Learning: Algorithm

Initialize q arbitrarily

Repeat (for each episode):

Initialize s

Repeat (for each time step of episode):

Take action from s according to policy (e.g. ϵ -greedy) $a \sim \pi(s)$

Observe reward r and transit to new state s'

$$q(s, a) \leftarrow (1 - \alpha)q(s, a) + \alpha \left(r + \gamma \max_{a'} q(s', a') \right)$$

$$s \leftarrow s'$$

until s is terminal



Q-Learning: Review

◎ Strength

- ◆ Learn from experience sampled from policy
- ◆ Choose the best action for each state to get highest long-term reward

◎ Weakness

- ◆ Over-estimate the action values using the maximum value as approximation for the maximum expected value
 - Double Q-learning
- ◆ Large memory to store q values with increasing numbers of states and actions
 - Deep Q-learning



Double Q-learning

◎ Double estimator

- ◆ sometimes underestimates rather than overestimates the maximum expected value

```
1: Initialize  $Q^A, Q^B, s$ 
2: repeat
3:   Choose  $a$ , based on  $Q^A(s, \cdot)$  and  $Q^B(s, \cdot)$ , observe  $r, s'$ 
4:   Choose (e.g. random) either UPDATE(A) or UPDATE(B)
5:   if UPDATE(A) then
6:     Define  $a^* = \arg \max_a Q^A(s', a)$ 
7:      $Q^A(s, a) \leftarrow Q^A(s, a) + \alpha(s, a) (r + \gamma Q^B(s', a^*) - Q^A(s, a))$ 
8:   else if UPDATE(B) then
9:     Define  $b^* = \arg \max_a Q^B(s', a)$ 
10:     $Q^B(s, a) \leftarrow Q^B(s, a) + \alpha(s, a) (r + \gamma Q^A(s', b^*) - Q^B(s, a))$ 
11:  end if
12:   $s \leftarrow s'$ 
13: until end
```



Function Approximation

- Represent value function by Q-network with weights θ

$$q(s, a, \theta) \approx q_*(s, a)$$

- May use back-propagation to compute $\nabla_{\theta} q$ in neural networks
- Popular choice: linear approximation, kernel methods, decision trees, **deep learning models**



Deep Q-Learning

◎ Advantage

- ◆ possible to apply the algorithm to larger problems, even when the state space is continuous
- ◆ generalize earlier experiences to previously unseen states

◎ Challenge

- ◆ unstable or divergent when a nonlinear function approximator such as a neural network is used to represent Q-value



Q-learning: Variants

- ◎ Deep Q Network (DQN)

- ◆ [Playing Atari with Deep Reinforcement Learning](#), V. Mnih et al., NIPS 2013 Workshop

- ◎ Double DQN

- ◆ [Deep Reinforcement Learning with Double Q-learning](#), H. van Hasselt et al., AAAI 2016

- ◎ Dueling DQN

- ◆ [Dueling Network Architectures for Deep Reinforcement Learning](#), Z. Wang et al., ICML 2016



Policy-Based RL

Value Based

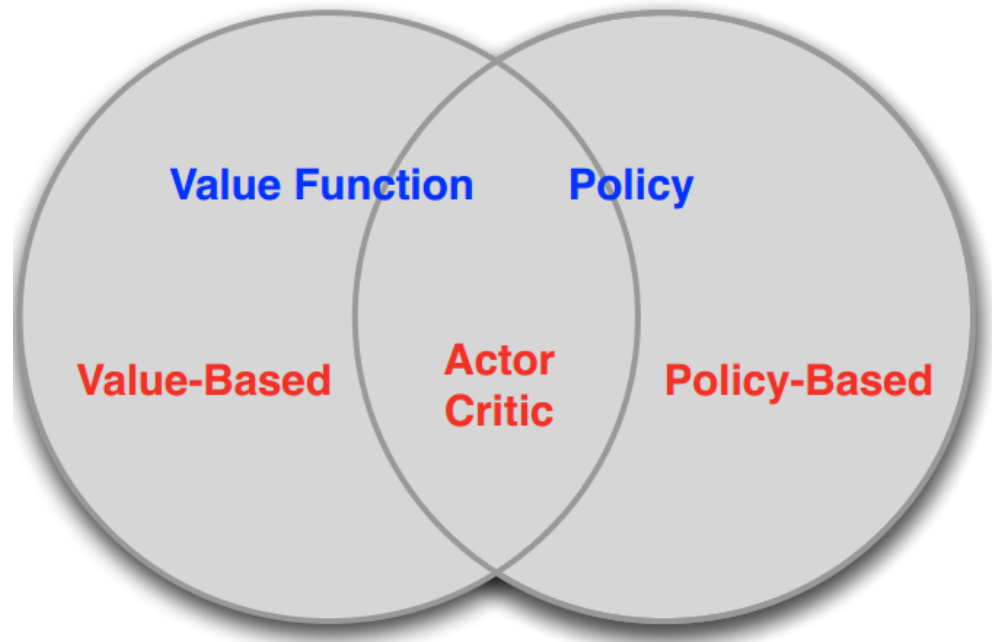
- ◆ Learnt Value Function
- ◆ Implicit Policy

Policy Based

- ◆ No Value Function
- ◆ Learnt Policy

Actor-Critic

- ◆ Learnt Value Function
- ◆ Learnt Policy



Policy-Based RL

⊙ Directly parametrize the **policy**: $\pi_{\theta}(s, a) = P[a|s, \theta]$

⊙ Comparing to Value-Based RL:

◆ Approximate the value function: $q_{\theta} \approx q^{\pi}(s, a)$

◆ Generate policy from the value function

- Greedy: taking the action that can maximize $q(s, a)$
- Epsilon-greedy: small probability to explore new

$$- a_t = \begin{cases} \arg \max_a q(s_t, a) & \text{with probability } 1 - \epsilon \\ \text{random action} & \text{with probability } \epsilon \end{cases}$$



Advantages

⊙ Advantages:

- ◆ Better convergence properties
- ◆ Effective in high-dimensional or continuous action spaces
- ◆ Can learn stochastic policies

⊙ Disadvantages:

- ◆ Converge to a local optimum
- ◆ High variance



Policy-based RL

- Markov Chain:

$$p_{\theta}(s_1, a_1, \dots, s_T, a_T) = p(s_1) \prod_{t=1}^T \pi_{\theta}(a_t | s_t) p(s_{t+1} | s_t, a_t)$$

- Goal is to maximize the reward:

$$\theta^* = \operatorname{argmax}_{\theta} E_{(s,a) \sim p_{\theta}(s,a)} [r(s, a)]$$

- Let $\tau = (s, a)$ be the state-action sequence:

$$\theta^* = \operatorname{argmax}_{\theta} E_{(\tau) \sim p_{\theta}(\tau)} [r(\tau)]$$



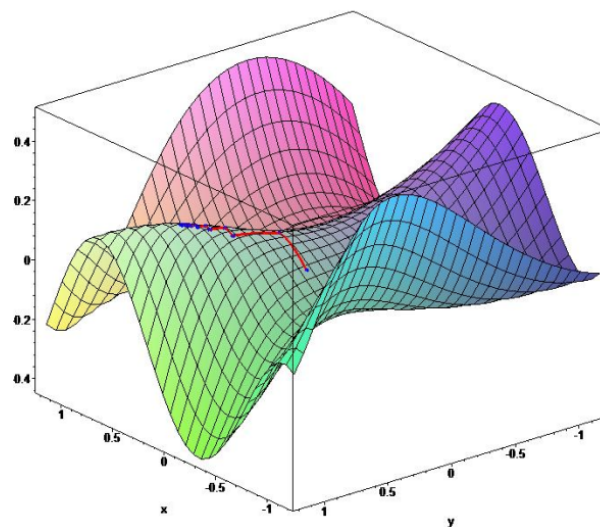
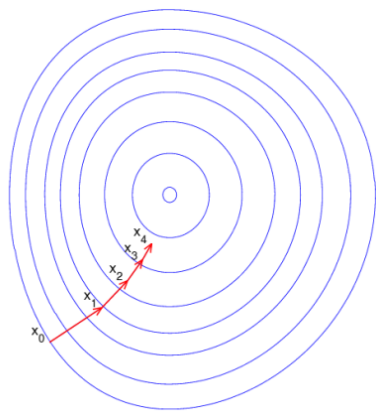
Policy Search

- Goal: find best θ for policy $\pi_{\theta}(\tau)$
- How to measure the quality of policy $\pi_{\theta}(\tau)$?
- Objective Function $J(\theta)$
 - In Policy Gradient, $J(\theta) = E_{\tau \sim \pi_{\theta}(\tau)} [\sum r(\tau)]$
- An **optimization** problem: find θ to maximize $J(\theta)$
 - Gradient Descent



Policy Gradient

- Search for a local maximum $J(\theta)$ by ascending the gradient of the policy: $\Delta\theta = \alpha \nabla_{\theta} J(\theta)$
- $\nabla_{\theta} J(\theta)$ is the **policy gradient**
- And α is learning rate



Policy Gradient

$$J(\theta) = E_{\tau \sim \pi_{\theta}(\tau)}[r(\tau)] = \int \pi_{\theta}(\tau) r(\tau) d\tau$$

$$\nabla_{\theta} J(\theta) = \int \nabla_{\theta} \pi_{\theta}(\tau) r(\tau) d\tau$$

$$= \int \pi_{\theta}(\tau) \nabla_{\theta} \log \pi_{\theta}(\tau) r(\tau) d\tau$$

$$\nabla_{\theta} \pi_{\theta}(\tau) = \pi_{\theta}(\tau) \frac{\nabla_{\theta} \pi_{\theta}(\tau)}{\pi_{\theta}(\tau)} = \pi_{\theta}(\tau) \nabla_{\theta} \log \pi_{\theta}(\tau)$$



Policy Gradient

$$\begin{aligned}\nabla_{\theta} J(\theta) &= \int \nabla_{\theta} \pi_{\theta}(\tau) r(\tau) d\tau \\&= \int \pi_{\theta}(\tau) \nabla_{\theta} \log \pi_{\theta}(\tau) r(\tau) d\tau \\&= E_{\tau \sim \pi_{\theta}(\tau)} [\nabla_{\theta} \log \pi_{\theta}(\tau) r(\tau)] \\&= E_{\tau \sim \pi_{\theta}(\tau)} \left[\left(\sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \right) \left(\sum_{t=1}^T r(s_t, a_t) \right) \right]\end{aligned}$$



One-Step MDPs

- ⊙ A simple class of one-step MDPs

- ◆ Start in state $\mathbf{s} \sim \mathbf{d}(\mathbf{s})$
- ◆ Obtain reward $\mathbf{r} = \mathbf{R}_{\mathbf{s},\mathbf{a}}$ after one time-step
- ◆ Then terminate

- ⊙ Compute **policy gradient** with likelihood ratios:

$$J(\theta) = E_{\tau \sim \pi_{\theta}(\tau)} \left[\sum r(\tau) \right] = \sum_{s \in S} d(s) \sum_{a \in A} \pi_{\theta}(\tau) R_{s,a}$$

$$\begin{aligned} \nabla_{\theta} J(\theta) &= \sum_{s \in S} d(s) \sum_{a \in A} \pi_{\theta}(\tau) \nabla_{\theta} \log \pi_{\theta}(\tau) R_{s,a} \\ &= E_{\tau \sim \pi_{\theta}(\tau)} [\nabla_{\theta} \log \pi_{\theta}(\tau) r(\tau)] \end{aligned}$$



Policy Gradient Theorem

- ◎ The policy gradient theorem:
 - ◆ Generalizes the likelihood ratio approach to multi-step MDPs
 - ◆ Replaces instantaneous reward r with long-term value $R^\pi(\tau)$

Policy Gradient Theorem

For any differentiable **policy** $\pi_\theta(\tau)$

For any **policy objective function** $J(\theta)$

the **policy gradient** is

$$\nabla_\theta J(\theta) = E_{\tau \sim \pi_\theta(\tau)} [\nabla_\theta \log \pi_\theta(\tau) R^\pi(\tau)]$$



REINFORCE

- How to update the policy with policy gradient?
- Using stochastic gradient ascent
- Using delayed reward v_t as an unbiased sample of $R^\pi(s, a)$

function REINFORCE

Initialise θ arbitrarily

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

$$\theta = \theta + \alpha \nabla_\theta \log \pi_\theta(s_t, a_t) v_t$$

end for

end for

Return θ

Review

- Policy-Based RL: directly parametrize the policy
- Objective Function $J(\theta)$: measure the quality of policy

$$J(\theta) = E_{\tau \sim \pi_{\theta}(\tau)} \left[\sum r(\tau) \right]$$

- Policy Gradient Theorem:

$$\nabla_{\theta} J(\theta) = E_{\tau \sim \pi_{\theta}(\tau)} [\nabla_{\theta} \log \pi_{\theta}(\tau) R^{\pi}(\tau)]$$

- REINFORCE algorithm: stochastic gradient ascent and delayed reward



Actor-Critic

Value Based

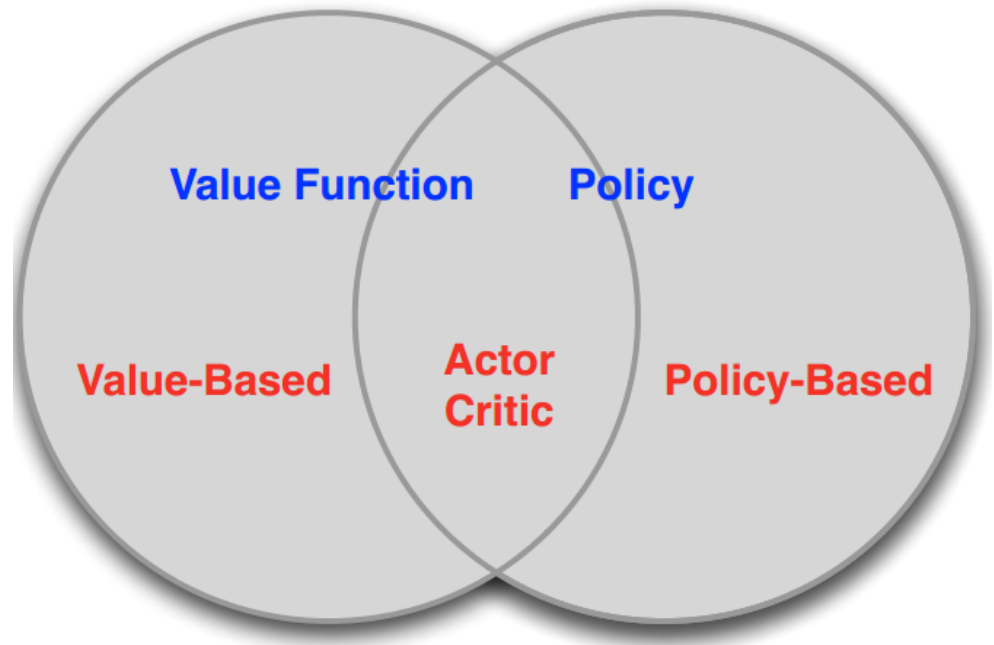
- ◆ Learnt Value Function
- ◆ Implicit Policy

Policy Based

- ◆ No Value Function
- ◆ Learnt Policy

Actor-Critic

- ◆ Learnt Value Function
- ◆ Learnt Policy



Actor-Critic

- Policy gradient still has high variance ☹
- Use a **critic** to evaluate the action just selected
- State-value function & state-action value function

$$V_w(s) = E_{a \sim \pi_\theta(a|s)}[q(a, s)]$$

- Actor-Critic maintains two sets of parameters
 - ◆ Actor (Policy Gradient) : Updates policy parameters θ
 - ◆ Critic (Function Approximation) : Updates parameters w of the state-value function



Actor-Critic

- Approximate policy gradient

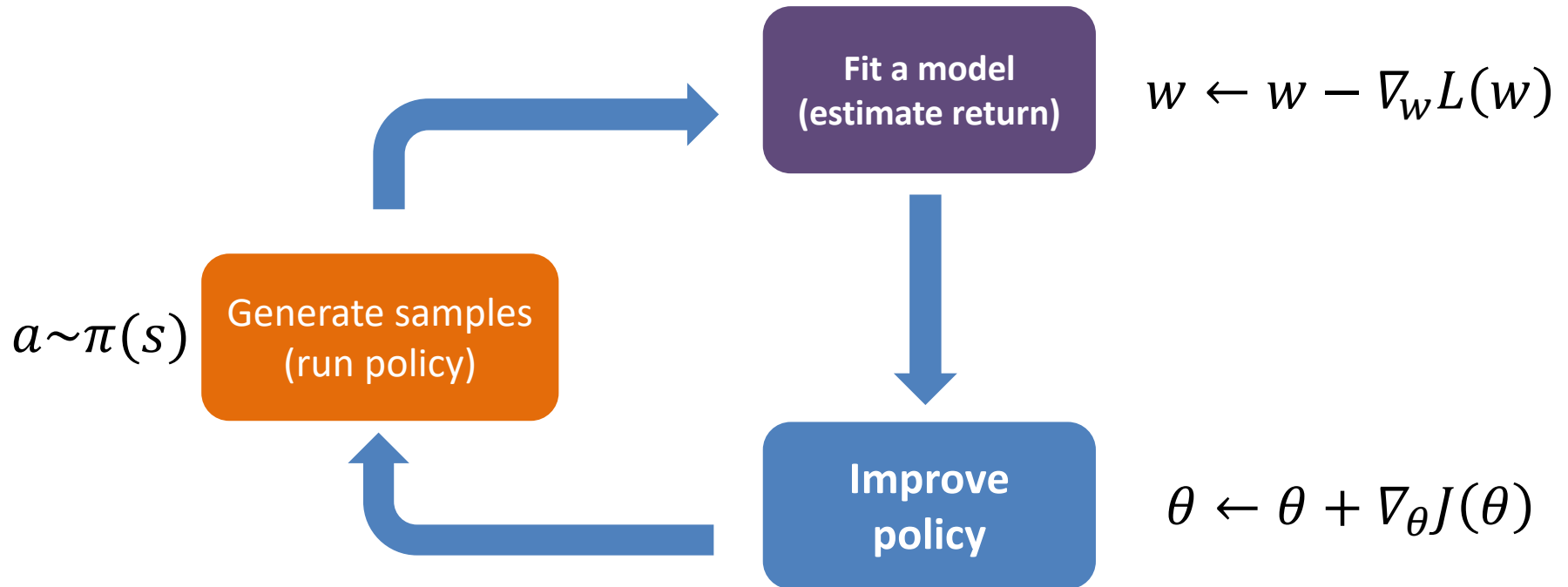
$$\begin{aligned}\nabla_{\theta} J(\theta) &= E_{\tau \sim \pi_{\theta}(\tau)} [\nabla_{\theta} \log \pi_{\theta}(\tau) R^{\pi}(\tau)] \\ &\rightarrow E_{\tau \sim \pi_{\theta}(\tau)} [\nabla_{\theta} \log \pi_{\theta}(\tau) \underbrace{(r + \gamma V_w(\tau') - V_w(\tau))}_{\text{TD error in Q-learning}}]\end{aligned}$$

- Fit $V_w(s)$ to sampled reward sums $\sum R^{\pi}(s', a') \approx R^{\pi}(s, a) + V_w(s')$

$$L(w) = \frac{1}{2} || \overbrace{r + \gamma V_w(s')}^{\text{Target}} - V_w(s) ||^2$$



Actor-Critic: Procedure



Actor-Critic: Algorithm

Initialize actor π with θ , critic V with w arbitrarily

Repeat (for each episode):

Initialize s

Repeat (for each time step of episode):

Take action from s according to actor $a \sim \pi_{\theta}(s)$

Observe reward r and transit to new state s'

Evaluate TD error $A(s) = r + \gamma V_w(s') - V_w(s)$

$w \leftarrow w - \alpha_w \nabla_w ||A(s)||^2$

Calculate $A(s)$ using new critic V_w

$\theta \leftarrow \theta + \alpha_{\theta} \nabla_{\theta} \log \pi_{\theta}(s) A(s)$

$s \leftarrow s'$

until s is terminal



◎ Summary and Discussions



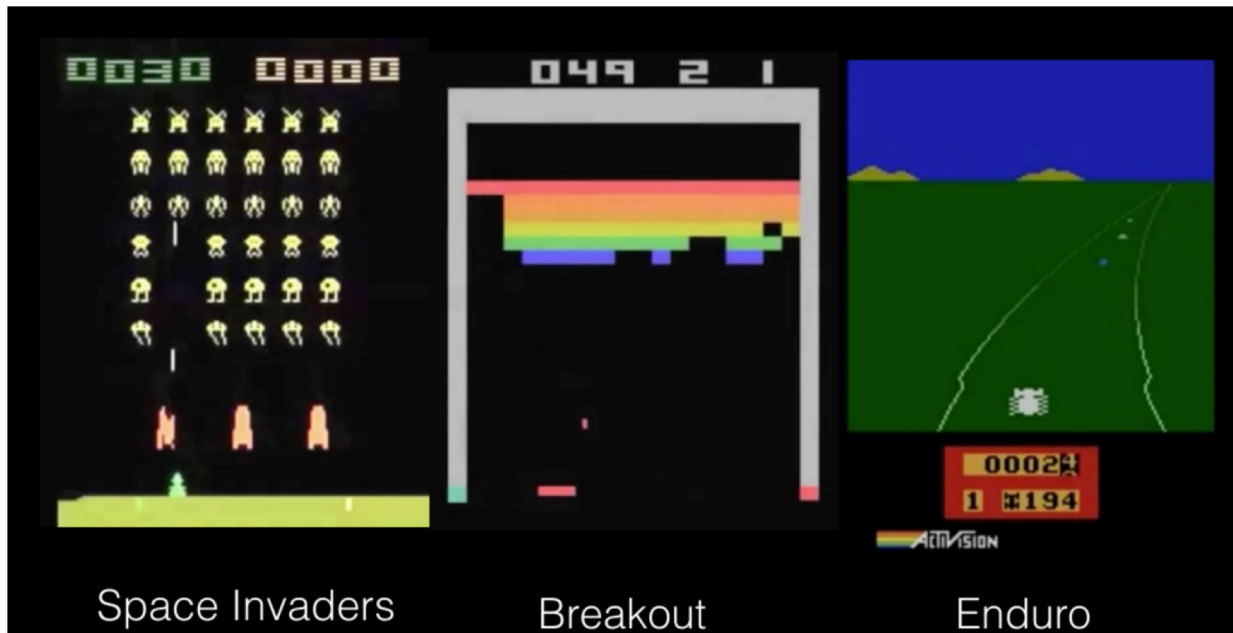
Reinforcement Learning

- ◎ **Sequential decision**: current decision affects future decision
- ◎ **Trial-and-error**: just try, do not worry making mistakes
 - ◆ **Explore** (new possibilities)
 - ◆ **Exploit** (with the current best policy)
- ◎ **Future reward**: maximizing the future rewards instead of just the intermediate rewards at each step
 - ◆ Remember $q(s,a)$



Difference to Supervised Learning

- Supervised learning: given a set of samples (x_i, y_i) , estimate $f: X \rightarrow Y$



Difference to Supervised Learning

- ⊙ You know what a true goal is, but do not know how to achieve that goal
- ⊙ Through interactions with environment (trial-and-error)
- ⊙ Many possible solutions (policies), which is optimal?

