

强化学习与博弈论

Reinforcement Learning and Game Theory

陈旭

计算机学院



中山大學
SUN YAT-SEN UNIVERSITY

Policy-Based Reinforcement Learning

Policy Function Approximation

Policy Function $\pi(a|s)$

- Policy function $\pi(a|s)$ is a probability density function (PDF).
- It takes state s as input.
- It output the probabilities for all the actions, e.g.,

$$\pi(\text{left}|s) = 0.2,$$

$$\pi(\text{right}|s) = 0.1,$$

$$\pi(\text{up}|s) = 0.7.$$

- Randomly sample action a random drawn from the distribution.

Can we directly learn a policy function $\pi(\textcolor{red}{a}|\textcolor{green}{s})$?

- If there are only a few states and actions, then yes, we can.
- Draw a table (matrix) and learn the entries.

| | Action a_1 | Action a_2 | Action a_3 | Action a_4 | ... |
|-------------|--------------|--------------|--------------|--------------|-----|
| State s_1 | | | | | |
| State s_2 | | | | | |
| State s_3 | | | | | |
| ⋮ | | | | | |

Can we directly learn a policy function $\pi(\textcolor{red}{a}|\textcolor{green}{s})$?

- If there are only a few states and actions, then yes, we can.
- Draw a table (matrix) and learn the entries.
- What if there are too many (or infinite) states or actions?

| | Action a_1 | Action a_2 | Action a_3 | Action a_4 | ... |
|-------------|--------------|--------------|--------------|--------------|-----|
| State s_1 | | | | | |
| State s_2 | | | | | |
| State s_3 | | | | | |
| ⋮ | | | | | |

Policy Network $\pi(a|s; \theta)$

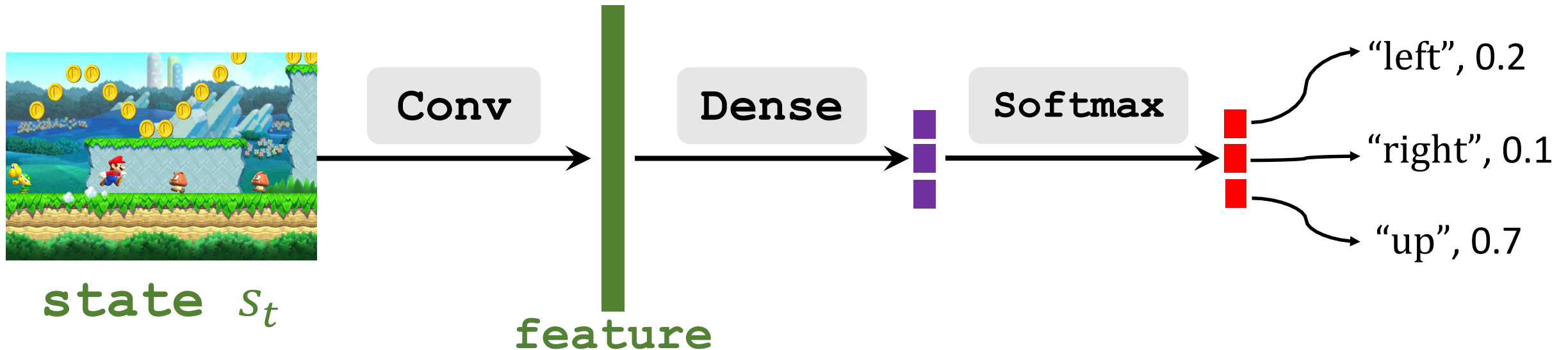
Policy network: Use a neural net to approximate $\pi(a|s)$.

- Use policy network $\pi(a|s; \theta)$ to approximate $\pi(a|s)$.
- θ : trainable parameters of the neural net.

Policy Network $\pi(a|s; \theta)$

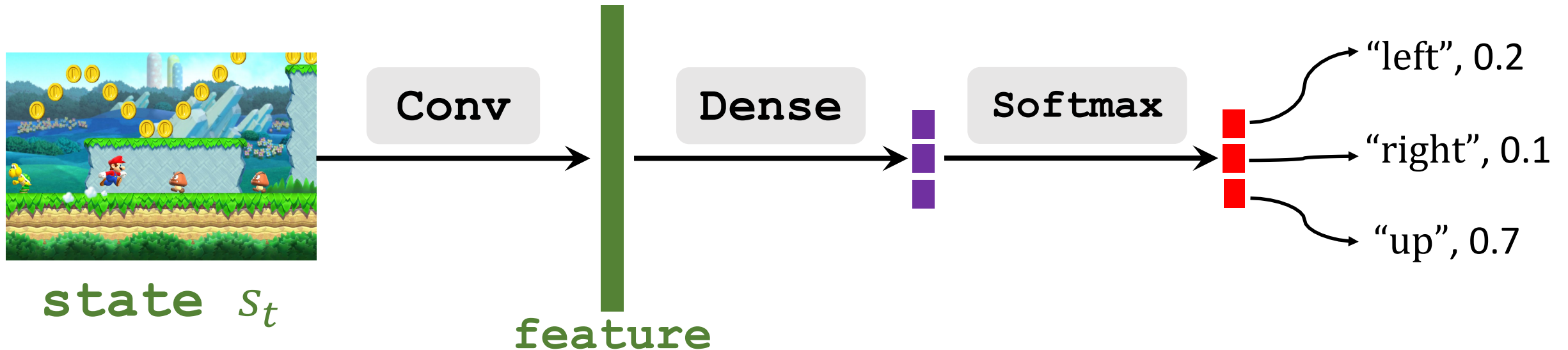
Policy network: Use a neural net to approximate $\pi(a|s)$.

- Use policy network $\pi(a|s; \theta)$ to approximate $\pi(a|s)$.
- θ : trainable parameters of the neural net.



Policy Network $\pi(a|s; \theta)$

- $\sum_{a \in \mathcal{A}} \pi(a|s; \theta) = 1$.
- Here, $\mathcal{A} = \{\text{"left"}, \text{"right"}, \text{"up"}\}$ is the set all actions.
- That is why we use softmax activation.



State-Value Function Approximation

Action-Value Function

Definition: Discounted return.

- $U_t = R_t + \gamma \cdot R_{t+1} + \gamma^2 \cdot R_{t+2} + \gamma^3 \cdot R_{t+3} + \dots$



- The return depends on actions $A_t, A_{t+1}, A_{t+2}, \dots$ and states $S_t, S_{t+1}, S_{t+2}, \dots$
- Actions are random: $\mathbb{P}[A = a \mid S = s] = \pi(a|s)$. (Policy function.)
- States are random: $\mathbb{P}[S' = s' \mid S = s, A = a] = p(s'|s, a)$. (State transition.)


Action-Value Function

Definition: Discounted return.

- $U_t = R_t + \gamma \cdot R_{t+1} + \gamma^2 \cdot R_{t+2} + \gamma^3 \cdot R_{t+3} + \dots$

Definition: Action-value function.

- $Q_\pi(s_t, a_t) = \mathbb{E} [U_t | S_t = s_t, A_t = a_t].$



The expectation is taken w.r.t.
actions $A_{t+1}, A_{t+2}, A_{t+3}, \dots$
and states $S_{t+1}, S_{t+2}, S_{t+3}, \dots$

State-Value Function

Definition: Discounted return.

- $U_t = R_t + \gamma \cdot R_{t+1} + \gamma^2 \cdot R_{t+2} + \gamma^3 \cdot R_{t+3} + \dots$

Definition: Action-value function.

- $Q_\pi(s_t, a_t) = \mathbb{E} [U_t | S_t = s_t, A_t = a_t].$

Definition: State-value function.

- $V_\pi(s_t) = \mathbb{E}_A [Q_\pi(s_t, A)]$



Integrate out action $A \sim \pi(\cdot | s_t).$

State-Value Function

Definition: Discounted return.

- $U_t = R_t + \gamma \cdot R_{t+1} + \gamma^2 \cdot R_{t+2} + \gamma^3 \cdot R_{t+3} + \dots$

Definition: Action-value function.

- $Q_\pi(s_t, a_t) = \mathbb{E} [U_t | S_t = s_t, A_t = a_t].$

Definition: State-value function.

- $V_\pi(s_t) = \mathbb{E}_A [Q_\pi(s_t, A)] = \sum_a \pi(a | s_t) \cdot Q_\pi(s_t, a).$



Integrate out action $A \sim \pi(\cdot | s_t).$

Policy-Based Reinforcement Learning

Definition: State-value function.

- $V_{\pi}(s_t) = \mathbb{E}_A [Q_{\pi}(s_t, A)] = \sum_a \pi(a|s_t) \cdot Q_{\pi}(s_t, a).$

Policy-Based Reinforcement Learning

Definition: State-value function.

- $V_{\pi}(s_t) = \mathbb{E}_A [Q_{\pi}(s_t, A)] = \sum_a \pi(a|s_t) \cdot Q_{\pi}(s_t, a).$

Approximate state-value function.

- Approximate policy function $\pi(a|s_t)$ by policy network $\pi(a|s_t; \theta).$

Policy-Based Reinforcement Learning

Definition: State-value function.

- $V_{\pi}(s_t) = \mathbb{E}_A [Q_{\pi}(s_t, A)] = \sum_a \pi(a|s_t) \cdot Q_{\pi}(s_t, a).$

Approximate state-value function.

- Approximate policy function $\pi(a|s_t)$ by policy network $\pi(a|s_t; \theta).$
- Approximate value function $V_{\pi}(s_t)$ by:

$$V(s_t; \theta) = \sum_a \pi(a|s_t; \theta) \cdot Q_{\pi}(s_t, a).$$

Policy-Based Reinforcement Learning

Definition: Approximate state-value function.

- $V(s; \theta) = \sum_a \pi(a|s; \theta) \cdot Q_\pi(s, a).$

Policy-based learning: Learn θ that maximizes $J(\theta) = \mathbb{E}_s[V(S; \theta)].$

Policy-Based Reinforcement Learning

Definition: Approximate state-value function.

- $V(s; \theta) = \sum_a \pi(a|s; \theta) \cdot Q_\pi(s, a).$

Policy-based learning: Learn θ that maximizes $J(\theta) = \mathbb{E}_s[V(S; \theta)].$

How to improve θ ? Policy gradient ascent!

- Observe state s .

- Update policy by: $\theta \leftarrow \theta + \beta \cdot \frac{\partial V(s; \theta)}{\partial \theta}$

Policy gradient

Policy Gradient

Reference

- Sutton and others: [Policy gradient methods for reinforcement learning with function approximation](#). In *NIPS*, 2000.

Policy Gradient

Definition: Approximate state-value function.

- $V(s; \theta) = \sum_a \pi(a|s; \theta) \cdot Q_\pi(s, a).$

Policy gradient: Derivative of $V(s; \theta)$ w.r.t. θ .

- $\frac{\partial V(s; \theta)}{\partial \theta}$

Policy Gradient

Definition: Approximate state-value function.

- $V(s; \theta) = \sum_a \pi(a|s; \theta) \cdot Q_\pi(s, a).$

Policy gradient: Derivative of $V(s; \theta)$ w.r.t. θ .

- $$\frac{\partial V(s; \theta)}{\partial \theta} = \frac{\partial \sum_a \pi(a|s; \theta) \cdot Q_\pi(s, a)}{\partial \theta}$$

Policy Gradient

Definition: Approximate state-value function.

- $V(s; \theta) = \sum_a \pi(a|s; \theta) \cdot Q_\pi(s, a).$

Policy gradient: Derivative of $V(s; \theta)$ w.r.t. θ .

- $$\frac{\partial V(s; \theta)}{\partial \theta} = \frac{\partial \sum_a \pi(a|s; \theta) \cdot Q_\pi(s, a)}{\partial \theta}$$
$$= \sum_a \frac{\partial \pi(a|s; \theta) \cdot Q_\pi(s, a)}{\partial \theta}$$

← Push derivative inside the summation

Policy Gradient

Definition: Approximate state-value function.

- $V(s; \theta) = \sum_a \pi(a|s; \theta) \cdot Q_\pi(s, a).$

Policy gradient: Derivative of $V(s; \theta)$ w.r.t. θ .

- $$\begin{aligned} \frac{\partial V(s; \theta)}{\partial \theta} &= \frac{\partial \sum_a \pi(a|s; \theta) \cdot Q_\pi(s, a)}{\partial \theta} \\ &= \sum_a \frac{\partial \pi(a|s; \theta) \cdot Q_\pi(s, a)}{\partial \theta} \\ &= \sum_a \frac{\partial \pi(a|s; \theta)}{\partial \theta} \cdot Q_\pi(s, a) \end{aligned}$$

← Pretend Q_π is independent of θ .
(It may not be true.)

Policy Gradient

Definition: Approximate state-value function.

- $V(s; \theta) = \sum_a \pi(a|s; \theta) \cdot Q_\pi(s, a).$

Policy gradient: Derivative of $V(s; \theta)$ w.r.t. θ .

- $\frac{\partial V(s; \theta)}{\partial \theta} = \sum_a \frac{\partial \pi(a|s; \theta)}{\partial \theta} \cdot Q_\pi(s, a)$

Policy Gradient

Policy Gradient

Definition: Approximate state-value function.

- $V(\textcolor{green}{s}; \boldsymbol{\theta}) = \sum_{\textcolor{red}{a}} \pi(\textcolor{red}{a}|\textcolor{green}{s}; \boldsymbol{\theta}) \cdot Q_{\pi}(\textcolor{green}{s}, \textcolor{red}{a}).$

Policy gradient: Derivative of $V(\textcolor{green}{s}; \boldsymbol{\theta})$ w.r.t. $\boldsymbol{\theta}$.

- $$\begin{aligned} \frac{\partial V(\textcolor{green}{s}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} &= \sum_{\textcolor{red}{a}} \boxed{\frac{\partial \pi(\textcolor{red}{a}|\textcolor{green}{s}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}}} Q_{\pi}(\textcolor{green}{s}, \textcolor{red}{a}) \\ &= \sum_{\textcolor{red}{a}} \boxed{\pi(\textcolor{red}{a}|\textcolor{green}{s}; \boldsymbol{\theta}) \cdot \frac{\partial \log \pi(\textcolor{red}{a}|\textcolor{green}{s}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}}} \cdot Q_{\pi}(\textcolor{green}{s}, \textcolor{red}{a}) \end{aligned}$$

Policy Gradient

Definition: Approximate state-value function.

- $V(s; \theta) = \sum_a \pi(a|s; \theta) \cdot Q_\pi(s, a).$

Policy gradient: Derivative of $V(s; \theta)$ w.r.t. θ .

- $$\begin{aligned} \frac{\partial V(s; \theta)}{\partial \theta} &= \sum_a \frac{\partial \pi(a|s; \theta)}{\partial \theta} \cdot Q_\pi(s, a) \\ &= \sum_a \pi(a|s; \theta) \cdot \boxed{\frac{\partial \log \pi(a|s; \theta)}{\partial \theta}} \cdot Q_\pi(s, a) \end{aligned}$$

- Chain rule: $\frac{\partial \log[\pi(\theta)]}{\partial \theta} = \frac{1}{\pi(\theta)} \cdot \frac{\partial \pi(\theta)}{\partial \theta}.$

Policy Gradient

Definition: Approximate state-value function.

- $V(\textcolor{green}{s}; \boldsymbol{\theta}) = \sum_{\textcolor{red}{a}} \pi(\textcolor{red}{a}|\textcolor{green}{s}; \boldsymbol{\theta}) \cdot Q_{\pi}(\textcolor{green}{s}, \textcolor{red}{a}).$

Policy gradient: Derivative of $V(\textcolor{green}{s}; \boldsymbol{\theta})$ w.r.t. $\boldsymbol{\theta}$.

- $$\begin{aligned} \frac{\partial V(\textcolor{green}{s}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} &= \sum_{\textcolor{red}{a}} \frac{\partial \pi(\textcolor{red}{a}|\textcolor{green}{s}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \cdot Q_{\pi}(\textcolor{green}{s}, \textcolor{red}{a}) \\ &= \sum_{\textcolor{red}{a}} \pi(\textcolor{red}{a}|\textcolor{green}{s}; \boldsymbol{\theta}) \cdot \frac{\partial \log \pi(\textcolor{red}{a}|\textcolor{green}{s}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \cdot Q_{\pi}(\textcolor{green}{s}, \textcolor{red}{a}) \end{aligned}$$

- Chain rule: $\frac{\partial \log[\pi(\boldsymbol{\theta})]}{\partial \boldsymbol{\theta}} = \frac{1}{\pi(\boldsymbol{\theta})} \cdot \frac{\partial \pi(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}}.$
- $\rightarrow \pi(\boldsymbol{\theta}) \cdot \frac{\partial \log[\pi(\boldsymbol{\theta})]}{\partial \boldsymbol{\theta}} = \pi(\boldsymbol{\theta}) \cdot \frac{1}{\pi(\boldsymbol{\theta})} \cdot \frac{\partial \pi(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}}$

Policy Gradient

Definition: Approximate state-value function.

- $V(s; \theta) = \sum_a \pi(a|s; \theta) \cdot Q_\pi(s, a).$

Policy gradient: Derivative of $V(s; \theta)$ w.r.t. θ .

- $$\begin{aligned} \frac{\partial V(s; \theta)}{\partial \theta} &= \sum_a \frac{\partial \pi(a|s; \theta)}{\partial \theta} \cdot Q_\pi(s, a) \\ &= \sum_a \pi(a|s; \theta) \cdot \frac{\partial \log \pi(a|s; \theta)}{\partial \theta} \cdot Q_\pi(s, a) \end{aligned}$$

- Chain rule: $\frac{\partial \log[\pi(\theta)]}{\partial \theta} = \frac{1}{\pi(\theta)} \cdot \frac{\partial \pi(\theta)}{\partial \theta}.$
- $\rightarrow \pi(\theta) \cdot \frac{\partial \log[\pi(\theta)]}{\partial \theta} = \cancel{\pi(\theta)} \cdot \cancel{\frac{1}{\pi(\theta)}} \cdot \frac{\partial \pi(\theta)}{\partial \theta}.$

Policy Gradient

Definition: Approximate state-value function.

- $V(s; \theta) = \sum_a \pi(a|s; \theta) \cdot Q_\pi(s, a).$

Policy gradient: Derivative of $V(s; \theta)$ w.r.t. θ .

- $$\begin{aligned} \frac{\partial V(s; \theta)}{\partial \theta} &= \sum_a \frac{\partial \pi(a|s; \theta)}{\partial \theta} Q_\pi(s, a) \\ &= \sum_a \pi(a|s; \theta) \cdot \frac{\partial \log \pi(a|s; \theta)}{\partial \theta} \cdot Q_\pi(s, a) \end{aligned}$$

- Chain rule: $\frac{\partial \log[\pi(\theta)]}{\partial \theta} = \frac{1}{\pi(\theta)} \cdot \frac{\partial \pi(\theta)}{\partial \theta}.$
- $\rightarrow \pi(\theta) \cdot \frac{\partial \log[\pi(\theta)]}{\partial \theta} = \cancel{\pi(\theta)} \cdot \frac{1}{\cancel{\pi(\theta)}} \cdot \frac{\partial \pi(\theta)}{\partial \theta}.$

Policy Gradient

Definition: Approximate state-value function.

- $V(\mathbf{s}; \boldsymbol{\theta}) = \sum_{\mathbf{a}} \pi(\mathbf{a}|\mathbf{s}; \boldsymbol{\theta}) \cdot Q_{\pi}(\mathbf{s}, \mathbf{a}).$

Policy gradient: Derivative of $V(\mathbf{s}; \boldsymbol{\theta})$ w.r.t. $\boldsymbol{\theta}$.

- $$\begin{aligned} \frac{\partial V(\mathbf{s}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} &= \sum_{\mathbf{a}} \frac{\partial \pi(\mathbf{a}|\mathbf{s}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \cdot Q_{\pi}(\mathbf{s}, \mathbf{a}) \\ &= \sum_{\mathbf{a}} \pi(\mathbf{a}|\mathbf{s}; \boldsymbol{\theta}) \cdot \frac{\partial \log \pi(\mathbf{a}|\mathbf{s}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \cdot Q_{\pi}(\mathbf{s}, \mathbf{a}) \\ &= \mathbb{E}_{\mathbf{A}} \left[\frac{\partial \log \pi(\mathbf{A}|\mathbf{s}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \cdot Q_{\pi}(\mathbf{s}, \mathbf{A}) \right]. \end{aligned}$$

The expectation is taken w.r.t. the random variable $\mathbf{A} \sim \pi(\cdot | \mathbf{s}; \boldsymbol{\theta})$.

Policy Gradient

Policy gradient:

$$\frac{\partial V(\mathbf{s}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = \mathbb{E}_{\mathbf{A} \sim \pi(\cdot | \mathbf{s}; \boldsymbol{\theta})} \left[\frac{\partial \log \pi(\mathbf{A} | \mathbf{s}, \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \cdot Q_{\pi}(\mathbf{s}, \mathbf{A}) \right].$$

Calculate Policy Gradient

Policy Gradient: $\frac{\partial V(s; \theta)}{\partial \theta} = \mathbb{E}_{A \sim \pi(\cdot | s; \theta)} \left[\frac{\partial \log \pi(A | s, \theta)}{\partial \theta} \cdot Q_{\pi}(s, A) \right].$

Calculate Policy Gradient

Policy Gradient:
$$\frac{\partial V(\mathbf{s}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = \mathbb{E}_{\mathbf{A} \sim \pi(\cdot | \mathbf{s}; \boldsymbol{\theta})} \left[\frac{\partial \log \pi(\mathbf{A} | \mathbf{s}, \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \cdot Q_{\pi}(\mathbf{s}, \mathbf{A}) \right].$$

1. Randomly sample an action $\hat{\mathbf{a}}$ according to $\pi(\cdot | \mathbf{s}; \boldsymbol{\theta})$.

Calculate Policy Gradient

Policy Gradient: $\frac{\partial V(s; \theta)}{\partial \theta} = \mathbb{E}_{A \sim \pi(\cdot | s; \theta)} \left[\frac{\partial \log \pi(A | s, \theta)}{\partial \theta} \cdot Q_{\pi}(s, A) \right].$

1. Randomly sample an action \hat{a} according to $\pi(\cdot | s; \theta)$.

2. Calculate $\mathbf{g}(\hat{a}, \theta) = \frac{\partial \log \pi(\hat{a} | s; \theta)}{\partial \theta} \cdot Q_{\pi}(s, \hat{a}).$

- By the definition of \mathbf{g} , $\mathbb{E}_A[\mathbf{g}(A, \theta)] = \frac{\partial V(s; \theta)}{\partial \theta}.$
- $\mathbf{g}(\hat{a}, \theta)$ is an unbiased estimate of $\frac{\partial V(s; \theta)}{\partial \theta}.$

Calculate Policy Gradient

Policy Gradient: $\frac{\partial V(\mathbf{s}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = \mathbb{E}_{\mathbf{A} \sim \pi(\cdot | \mathbf{s}; \boldsymbol{\theta})} \left[\frac{\partial \log \pi(\mathbf{A} | \mathbf{s}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \cdot Q_{\pi}(\mathbf{s}, \mathbf{A}) \right].$

1. Randomly sample an action $\hat{\mathbf{a}}$ according to $\pi(\cdot | \mathbf{s}; \boldsymbol{\theta})$.
2. Calculate $\mathbf{g}(\hat{\mathbf{a}}, \boldsymbol{\theta}) = \frac{\partial \log \pi(\hat{\mathbf{a}} | \mathbf{s}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \cdot Q_{\pi}(\mathbf{s}, \hat{\mathbf{a}})$.
3. Use $\mathbf{g}(\hat{\mathbf{a}}, \boldsymbol{\theta})$ as an approximation to the policy gradient $\frac{\partial V(\mathbf{s}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}}$.

Update **policy network** using **policy gradient**

Algorithm

1. Observe the state s_t .
2. Randomly sample action a_t according to $\pi(\cdot | s_t; \theta_t)$.

Algorithm

1. Observe the state s_t .
2. Randomly sample action a_t according to $\pi(\cdot | s_t; \theta_t)$.
3. Compute $q_t \approx Q_\pi(s_t, a_t)$ (some estimate).
4. Differentiate policy network: $\mathbf{d}_{\theta,t} = \frac{\partial \log \pi(a_t | s_t, \theta)}{\partial \theta} \bigg|_{\theta=\theta_t}$.

Algorithm

1. Observe the state s_t .
2. Randomly sample action a_t according to $\pi(\cdot | s_t; \theta_t)$.
3. Compute $q_t \approx Q_\pi(s_t, a_t)$ (some estimate).
4. Differentiate policy network: $\mathbf{d}_{\theta,t} = \frac{\partial \log \pi(a_t | s_t, \theta)}{\partial \theta} \big|_{\theta=\theta_t}$.
5. (Approximate) policy gradient: $\mathbf{g}(a_t, \theta_t) = q_t \cdot \mathbf{d}_{\theta,t}$.
6. Update policy network: $\theta_{t+1} = \theta_t + \beta \cdot \mathbf{g}(a_t, \theta_t)$.

Algorithm

1. Observe the state s_t .
2. Randomly sample action a_t according to $\pi(\cdot | s_t; \theta_t)$.
3. Compute $q_t \approx Q_\pi(s_t, a_t)$ (some estimate). **How?**
4. Differentiate policy network: $\mathbf{d}_{\theta,t} = \frac{\partial \log \pi(a_t | s_t, \theta)}{\partial \theta} \big|_{\theta=\theta_t}$.
5. (Approximate) policy gradient: $\mathbf{g}(a_t, \theta_t) = q_t \cdot \mathbf{d}_{\theta,t}$.
6. Update policy network: $\theta_{t+1} = \theta_t + \beta \cdot \mathbf{g}(a_t, \theta_t)$.

Algorithm

1. Observe the state s_t .
2. Randomly sample action a_t according to $\pi(\cdot | s_t; \theta_t)$.
3. Compute $q_t \approx Q_\pi(s_t, a_t)$ (some estimate). **How?**

Option 1: REINFORCE.

- Play the game to the end and generate the trajectory:

$$s_1, a_1, r_1, s_2, a_2, r_2, \dots, s_T, a_T, r_T.$$

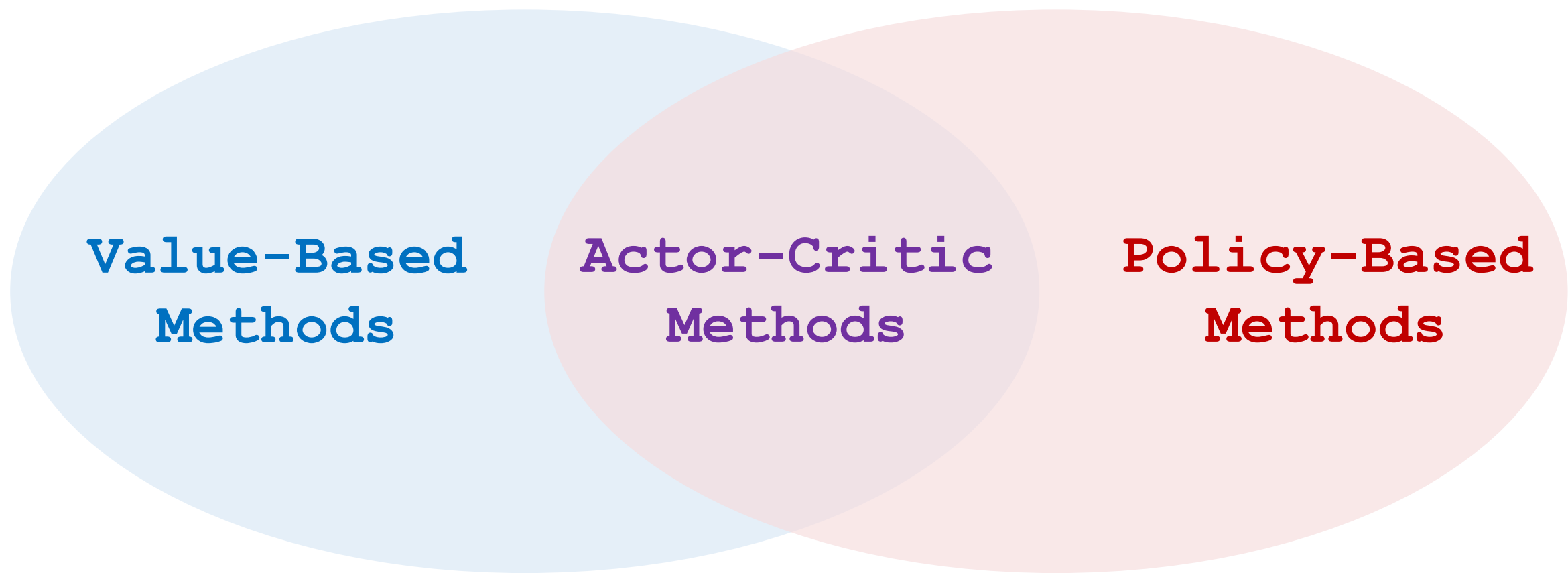
- Compute the discounted return $u_t = \sum_{k=t}^T \gamma^{k-t} r_k$, for all t .
- Since $Q_\pi(s_t, a_t) = \mathbb{E}[U_t]$, we can use u_t to approximate $Q_\pi(s_t, a_t)$.
- \Rightarrow Use $q_t = u_t$.

Algorithm

1. Observe the state s_t .
2. Randomly sample action a_t according to $\pi(\cdot | s_t; \theta_t)$.
3. Compute $q_t \approx Q_\pi(s_t, a_t)$ (some estimate). **How?**

Option 2: Approximate Q_π using a neural network.

- This leads to the actor-critic method.



Value Network and Policy Network

State-Value Function Approximation

Definition: State-value function.

- $V_{\pi}(s) = \sum_a \pi(a|s) \cdot Q_{\pi}(s, a).$

State-Value Function Approximation

Definition: State-value function.

- $V_{\pi}(s) = \sum_a \pi(a|s) \cdot Q_{\pi}(s, a).$

Policy network (actor):

- Use neural net $\pi(a|s; \theta)$ to approximate $\pi(a|s).$
- θ : trainable parameters of the neural net.

State-Value Function Approximation

Definition: State-value function.

- $V_{\pi}(s) = \sum_a \pi(a|s) \cdot Q_{\pi}(s, a).$

Policy network (actor):

- Use neural net $\pi(a|s; \theta)$ to approximate $\pi(a|s).$
- θ : trainable parameters of the neural net.

Value network (critic):

- Use neural net $q(s, a; \mathbf{w})$ to approximate $Q_{\pi}(s, a).$
- \mathbf{w} : trainable parameters of the neural net.

State-Value Function Approximation

Definition: State-value function.

- $V_{\pi}(s) = \sum_a \pi(a|s) \cdot Q_{\pi}(s, a) \approx \sum_a \pi(a|s; \theta) \cdot q(s, a; \mathbf{w}).$

Policy network (actor):

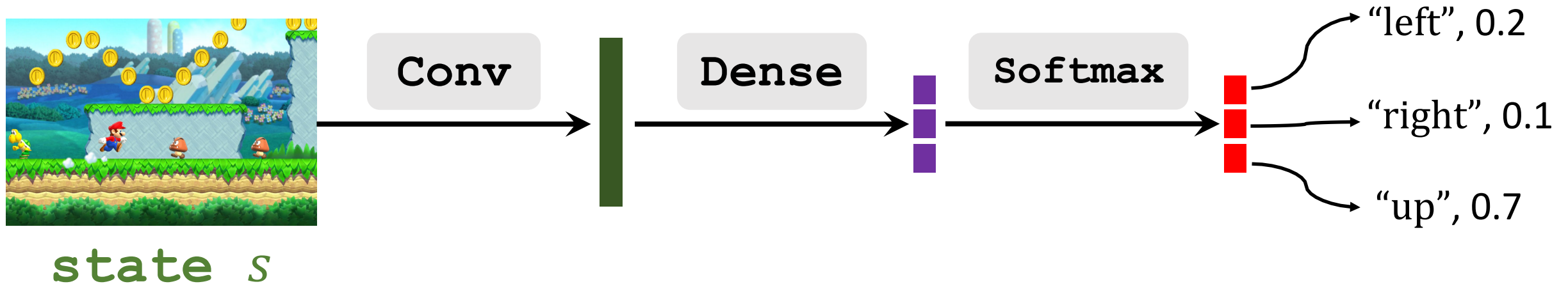
- Use neural net $\pi(a|s; \theta)$ to approximate $\pi(a|s)$.
- θ : trainable parameters of the neural net.

Value network (critic):

- Use neural net $q(s, a; \mathbf{w})$ to approximate $Q_{\pi}(s, a)$.
- \mathbf{w} : trainable parameters of the neural net.

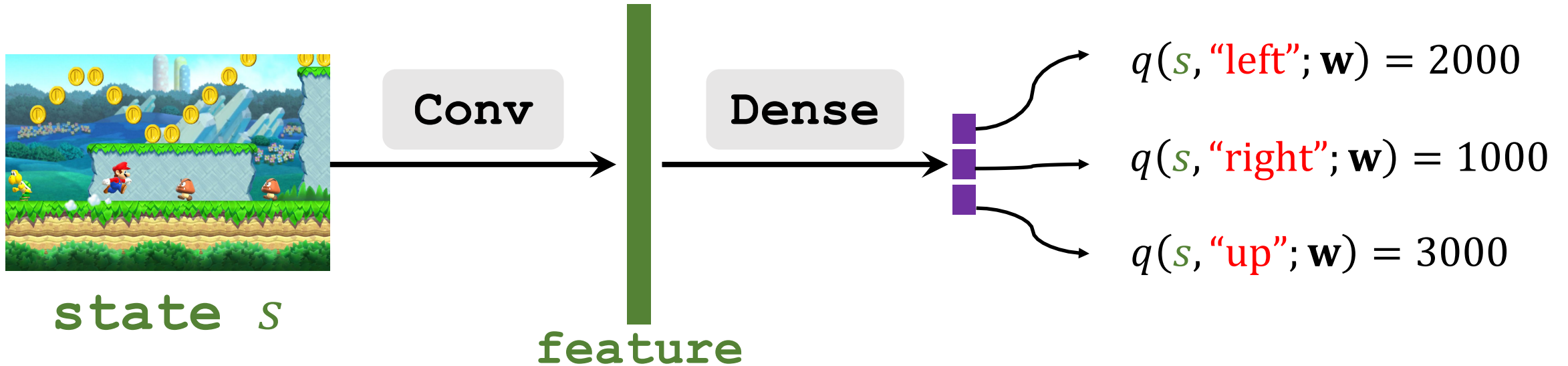
Policy Network (Actor): $\pi(a|s; \theta)$

- Input: **state** s , e.g., a screenshot of Super Mario.
- Output: probability distribution over the **actions**.
- Let \mathcal{A} be the set all actions, e.g., $\mathcal{A} = \{\text{"left"}, \text{"right"}, \text{"up"}\}$.
- $\sum_{a \in \mathcal{A}} \pi(a|s; \theta) = 1$. (That is why we use softmax activation.)



Value Network (Critic): $q(s, a; \mathbf{w})$

- Inputs: state s .
- Output: action-values of all the actions.



Actor-Critic Method

policy network (actor)



value network (critic)



Train the Neural Networks

Train the networks

Definition: State-value function approximated using neural networks.

- $V(\textcolor{green}{s}; \boldsymbol{\theta}, \mathbf{w}) = \sum_{\textcolor{red}{a}} \pi(\textcolor{red}{a} | \textcolor{green}{s}; \boldsymbol{\theta}) \cdot q(\textcolor{green}{s}, \textcolor{red}{a}; \mathbf{w})$.

Training: Update the parameters $\boldsymbol{\theta}$ and \mathbf{w} .

Train the networks

Definition: State-value function approximated using neural networks.

- $V(s; \theta, \mathbf{w}) = \sum_a \pi(a|s; \theta) \cdot q(s, a; \mathbf{w})$.

Training: Update the parameters θ and \mathbf{w} .

- Update policy network $\pi(a|s; \theta)$ to increase the state-value $V(s; \theta, \mathbf{w})$.
 - Actor gradually performs better.
 - Supervision is purely from the value network (critic).

Train the networks

Definition: State-value function approximated using neural networks.

- $V(s; \theta, \mathbf{w}) = \sum_a \pi(a|s; \theta) \cdot q(s, a; \mathbf{w})$.

Training: Update the parameters θ and \mathbf{w} .

- Update policy network $\pi(a|s; \theta)$ to increase the state-value $V(s; \theta, \mathbf{w})$.
 - Actor gradually performs better.
 - Supervision is purely from the value network (critic).
- Update value network $q(s, a; \mathbf{w})$ to better estimate the return.
 - Critic's judgement becomes more accurate.
 - Supervision is purely from the rewards.

Train the networks

Definition: State-value function approximated using neural networks.

- $V(\textcolor{green}{s}; \boldsymbol{\theta}, \mathbf{w}) = \sum_{\textcolor{red}{a}} \pi(\textcolor{red}{a} | \textcolor{green}{s}; \boldsymbol{\theta}) \cdot q(\textcolor{green}{s}, \textcolor{red}{a}; \mathbf{w})$.

Training: Update the parameters $\boldsymbol{\theta}$ and \mathbf{w} .

1. Observe the state $\textcolor{green}{s}_t$.
2. Randomly sample action $\textcolor{red}{a}_t$ according to $\pi(\cdot | \textcolor{green}{s}_t; \boldsymbol{\theta}_t)$.
3. Perform $\textcolor{red}{a}_t$ and observe new state $\textcolor{green}{s}_{t+1}$ and reward $\textcolor{blue}{r}_t$.
4. Update \mathbf{w} (in value network) using temporal difference (TD).
5. Update $\boldsymbol{\theta}$ (in policy network) using policy gradient.

Update **value network** q using **TD**

- Compute $q(s_t, a_t; \mathbf{w}_t)$ and $q(s_{t+1}, a_{t+1}; \mathbf{w}_t)$.
- TD target: $y_t = r_t + \gamma \cdot q(s_{t+1}, a_{t+1}; \mathbf{w}_t)$.

Update **value network** q using **TD**

- Compute $q(s_t, a_t; \mathbf{w}_t)$ and $q(s_{t+1}, a_{t+1}; \mathbf{w}_t)$.
- TD target: $y_t = r_t + \gamma \cdot q(s_{t+1}, a_{t+1}; \mathbf{w}_t)$.
- Loss: $L(\mathbf{w}) = \frac{1}{2} [q(s_t, a_t; \mathbf{w}) - y_t]^2$.
- Gradient descent: $\mathbf{w}_{t+1} = \mathbf{w}_t - \alpha \cdot \frac{\partial L(\mathbf{w})}{\partial \mathbf{w}} \Big|_{\mathbf{w}=\mathbf{w}_t}$.

Update **policy network** π using **policy gradient**

Definition: State-value function approximated using neural networks.

- $V(\textcolor{green}{s}; \boldsymbol{\theta}, \mathbf{w}) = \sum_{\textcolor{red}{a}} \pi(\textcolor{red}{a} | \textcolor{green}{s}; \boldsymbol{\theta}) \cdot q(\textcolor{green}{s}, \textcolor{red}{a}; \mathbf{w})$.

Policy gradient: Derivative of $V(\textcolor{green}{s}_t; \boldsymbol{\theta}, \mathbf{w})$ w.r.t. $\boldsymbol{\theta}$.

- Let $\mathbf{g}(\textcolor{red}{a}, \boldsymbol{\theta}) = \frac{\partial \log \pi(\textcolor{red}{a} | \textcolor{green}{s}, \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \cdot q(\textcolor{green}{s}_t, \textcolor{red}{a}; \mathbf{w})$.
- $\frac{\partial V(\textcolor{green}{s}; \boldsymbol{\theta}, \mathbf{w}_t)}{\partial \boldsymbol{\theta}} = \mathbb{E}_{\textcolor{red}{A}}[\mathbf{g}(\textcolor{red}{A}, \boldsymbol{\theta})]$.

Update **policy network** π using **policy gradient**

Definition: State-value function approximated using neural networks.

- $V(\textcolor{green}{s}; \boldsymbol{\theta}, \mathbf{w}) = \sum_{\textcolor{red}{a}} \pi(\textcolor{red}{a} | \textcolor{green}{s}; \boldsymbol{\theta}) \cdot q(\textcolor{green}{s}, \textcolor{red}{a}; \mathbf{w})$.

Policy gradient: Derivative of $V(\textcolor{green}{s}_t; \boldsymbol{\theta}, \mathbf{w})$ w.r.t. $\boldsymbol{\theta}$.

- Let $\mathbf{g}(\textcolor{red}{a}, \boldsymbol{\theta}) = \frac{\partial \log \pi(\textcolor{red}{a} | \textcolor{green}{s}, \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \cdot q(\textcolor{green}{s}_t, \textcolor{red}{a}; \mathbf{w})$.
- $\frac{\partial V(\textcolor{green}{s}; \boldsymbol{\theta}, \mathbf{w}_t)}{\partial \boldsymbol{\theta}} = \mathbb{E}_{\textcolor{red}{A}}[\mathbf{g}(\textcolor{red}{A}, \boldsymbol{\theta})]$.

Algorithm: Update policy network using stochastic policy gradient.

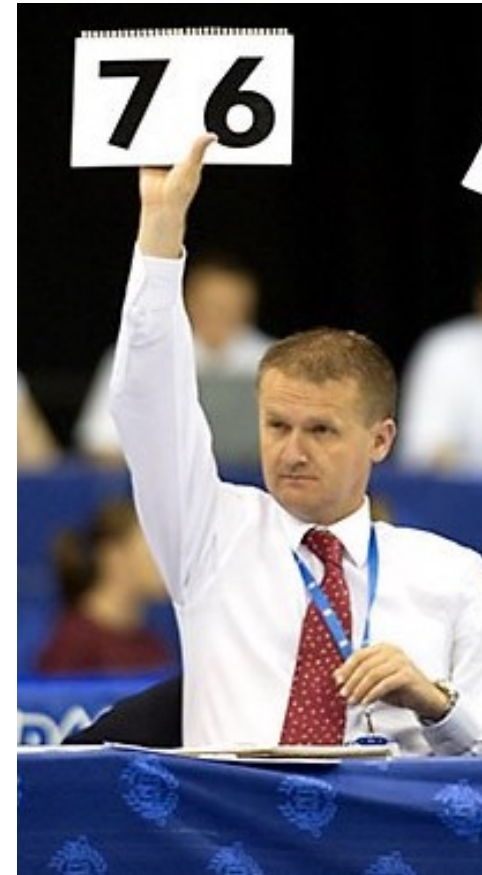
- Random sampling: $\textcolor{red}{a} \sim \pi(\cdot | \textcolor{green}{s}_t; \boldsymbol{\theta}_t)$. (Thus $\mathbf{g}(\textcolor{red}{a}, \boldsymbol{\theta})$ is unbiased.)
- Stochastic gradient ascent: $\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t + \beta \cdot \mathbf{g}(\textcolor{red}{a}, \boldsymbol{\theta}_t)$.

Actor-Critic Method

policy network (actor)



value network (critic)

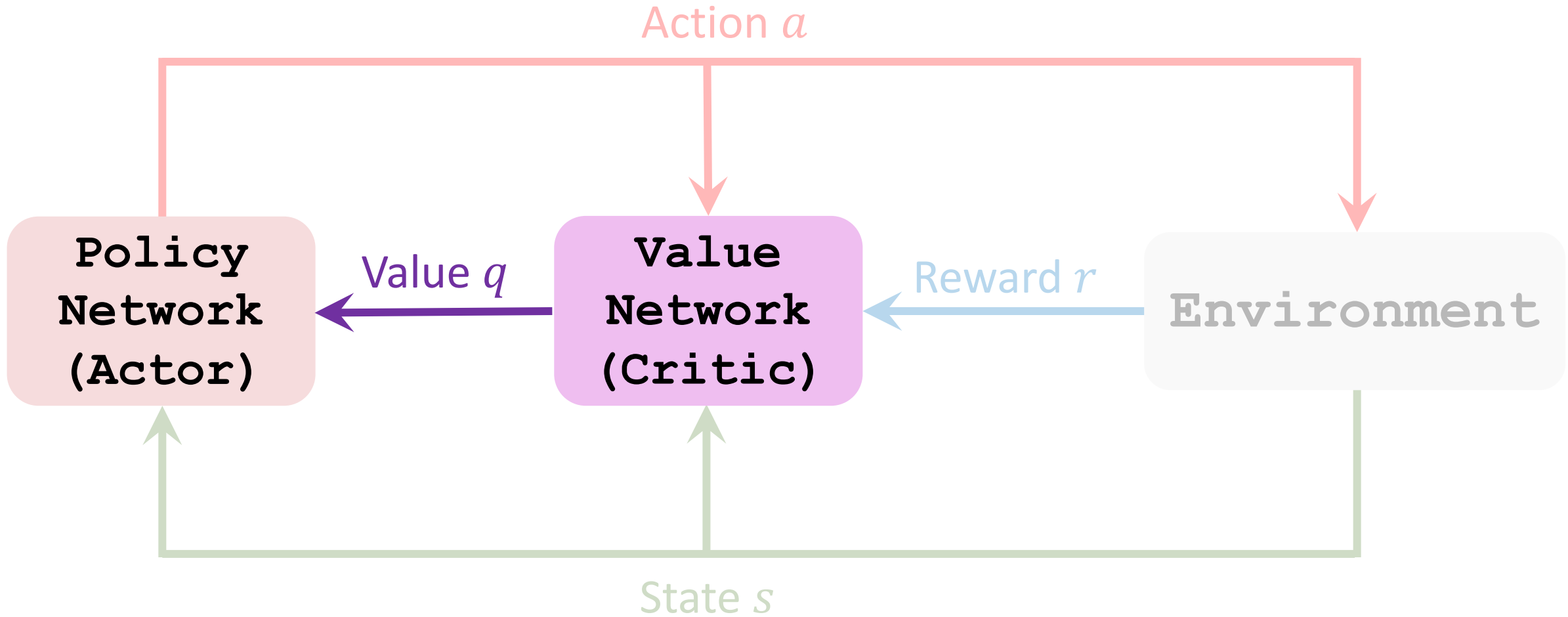


Actor-Critic Method

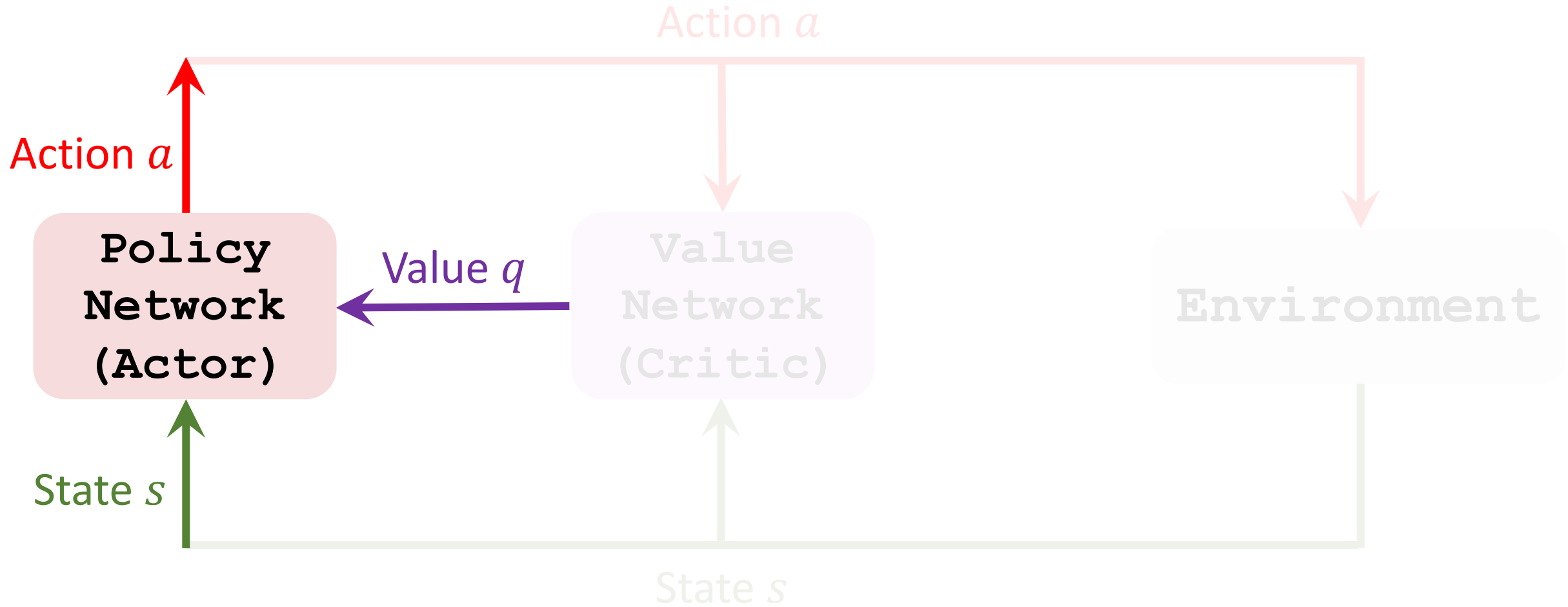
Action a



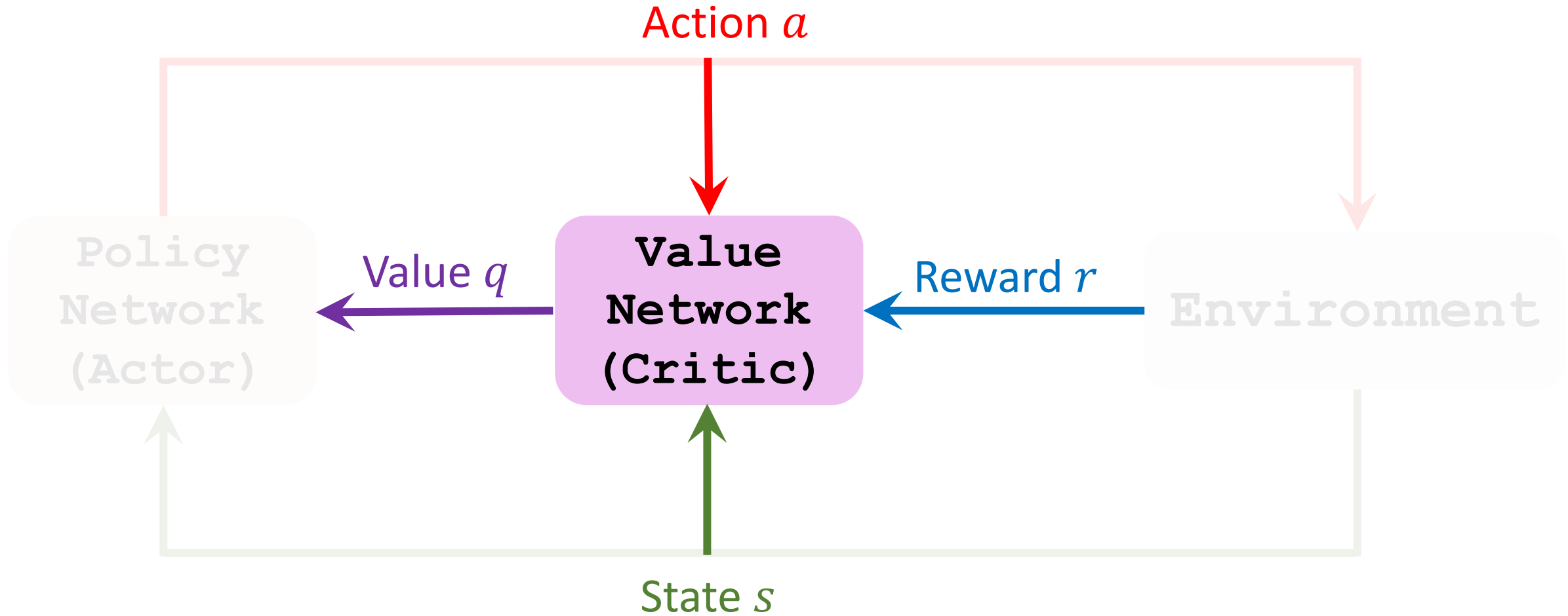
Actor-Critic Method



Actor-Critic Method: Update Actor



Actor-Critic Method: Update Critic



Summary of Algorithm

1. Observe state s_t and randomly sample $a_t \sim \pi(\cdot | s_t; \theta_t)$.
2. Perform a_t ; then environment gives new state s_{t+1} and reward r_t .
3. Randomly sample $\tilde{a}_{t+1} \sim \pi(\cdot | s_{t+1}; \theta_t)$. (Do not perform \tilde{a}_{t+1} !)

Summary of Algorithm

1. Observe state s_t and randomly sample $a_t \sim \pi(\cdot | s_t; \theta_t)$.
2. Perform a_t ; then environment gives new state s_{t+1} and reward r_t .
3. Randomly sample $\tilde{a}_{t+1} \sim \pi(\cdot | s_{t+1}; \theta_t)$. (Do not perform \tilde{a}_{t+1} !)
4. Evaluate value network: $q_t = q(s_t, a_t; \mathbf{w}_t)$ and $q_{t+1} = q(s_{t+1}, \tilde{a}_{t+1}; \mathbf{w}_t)$.
5. Compute TD error: $\delta_t = q_t - \underbrace{(r_t + \gamma \cdot q_{t+1})}_{\text{TD Target}}$.

Summary of Algorithm

1. Observe state s_t and randomly sample $a_t \sim \pi(\cdot | s_t; \theta_t)$.
2. Perform a_t ; then environment gives new state s_{t+1} and reward r_t .
3. Randomly sample $\tilde{a}_{t+1} \sim \pi(\cdot | s_{t+1}; \theta_t)$. (Do not perform \tilde{a}_{t+1} !)
4. Evaluate value network: $q_t = q(s_t, a_t; \mathbf{w}_t)$ and $q_{t+1} = q(s_{t+1}, \tilde{a}_{t+1}; \mathbf{w}_t)$.
5. Compute TD error: $\delta_t = q_t - (r_t + \gamma \cdot q_{t+1})$.
6. Differentiate value network: $\mathbf{d}_{w,t} = \frac{\partial q(s_t, a_t; \mathbf{w})}{\partial \mathbf{w}} \Big|_{\mathbf{w}=\mathbf{w}_t}$.
7. Update value network: $\mathbf{w}_{t+1} = \mathbf{w}_t - \alpha \cdot \delta_t \cdot \mathbf{d}_{w,t}$.

Summary of Algorithm

1. Observe state s_t and randomly sample $a_t \sim \pi(\cdot | s_t; \theta_t)$.
2. Perform a_t ; then environment gives new state s_{t+1} and reward r_t .
3. Randomly sample $\tilde{a}_{t+1} \sim \pi(\cdot | s_{t+1}; \theta_t)$. (Do not perform \tilde{a}_{t+1} !)
4. Evaluate value network: $q_t = q(s_t, a_t; \mathbf{w}_t)$ and $q_{t+1} = q(s_{t+1}, \tilde{a}_{t+1}; \mathbf{w}_t)$.
5. Compute TD error: $\delta_t = q_t - (r_t + \gamma \cdot q_{t+1})$.
6. Differentiate value network: $\mathbf{d}_{w,t} = \frac{\partial q(s_t, a_t; \mathbf{w})}{\partial \mathbf{w}} \Big|_{\mathbf{w}=\mathbf{w}_t}$.
7. Update value network: $\mathbf{w}_{t+1} = \mathbf{w}_t - \alpha \cdot \delta_t \cdot \mathbf{d}_{w,t}$.
8. Differentiate policy network: $\mathbf{d}_{\theta,t} = \frac{\partial \log \pi(a_t | s_t, \theta)}{\partial \theta} \Big|_{\theta=\theta_t}$.
9. Update policy network: $\theta_{t+1} = \theta_t + \beta \cdot q_t \cdot \mathbf{d}_{\theta,t}$.

Summary

Policy Network and Value Network

Definition: State-value function.

- $V_{\pi}(s) = \sum_a \pi(a|s) \cdot Q_{\pi}(s, a).$

Definition: function approximation using neural networks.

- Approximate policy function $\pi(a|s)$ by $\pi(a|s; \theta)$ (actor).
- Approximate value function $Q_{\pi}(s, a)$ by $q(s, a; \mathbf{w})$ (critic).

Roles of **Actor** and **Critic**

During training

- Agent is controlled by policy network (**actor**): $a_t \sim \pi(\cdot | s_t; \theta)$.
- Value network q (**critic**) provides the **actor** with supervision.

Roles of **Actor** and **Critic**

During training

- Agent is controlled by policy network (**actor**): $a_t \sim \pi(\cdot | s_t; \theta)$.
- Value network q (**critic**) provides the **actor** with supervision.

After training

- Agent is controlled by policy network (**actor**): $a_t \sim \pi(\cdot | s_t; \theta)$.
- Value network q (**critic**) will not be used.

Training

Update the **policy network (actor)** by **policy gradient**.

- Seek to increase state-value: $V(\mathbf{s}; \boldsymbol{\theta}, \mathbf{w}) = \sum_a \pi(a|\mathbf{s}; \boldsymbol{\theta}) \cdot q(\mathbf{s}, a; \mathbf{w})$.
- Compute policy gradient: $\frac{\partial V(\mathbf{s}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = \mathbb{E}_A \left[\frac{\partial \log \pi(A|\mathbf{s}, \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \cdot q(\mathbf{s}, A; \mathbf{w}) \right]$.
- Perform gradient ascent.

Training

Update the **policy network (actor)** by **policy gradient**.

- Seek to increase state-value: $V(\mathbf{s}; \boldsymbol{\theta}, \mathbf{w}) = \sum_a \pi(a|\mathbf{s}; \boldsymbol{\theta}) \cdot q(\mathbf{s}, a; \mathbf{w})$.
- Compute policy gradient: $\frac{\partial V(\mathbf{s}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = \mathbb{E}_A \left[\frac{\partial \log \pi(A|\mathbf{s}, \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \cdot q(\mathbf{s}, A; \mathbf{w}) \right]$.
- Perform gradient ascent.

Update the **value network (critic)** by **TD learning**.

- Predicted action-value: $q_t = q(\mathbf{s}_t, \mathbf{a}_t; \mathbf{w})$.
- TD target: $y_t = r_t + \gamma \cdot q(\mathbf{s}_{t+1}, \mathbf{a}_{t+1}; \mathbf{w})$
- Gradient: $\frac{\partial (q_t - y_t)^2 / 2}{\partial \mathbf{w}} = (q_t - y_t) \cdot \frac{\partial q(\mathbf{s}_t, \mathbf{a}_t; \mathbf{w})}{\partial \mathbf{w}}$.
- Perform gradient descent.