

Policy-Based Reinforcement Learning

Shusen Wang

Policy Function Approximation


Action-Value Function

Definition: Discounted return (aka cumulative discounted future reward).

- $R_t = r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \gamma^3 r_{t+3} + \dots$ (to infinity.)

Definition: Action-value function.

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



- Taken w.r.t. actions $a_{t+1}, a_{t+2}, a_{t+3}, \dots$ and states $s_{t+1}, s_{t+2}, s_{t+3}, \dots$
- Actions are randomly sampled: $a_t \sim \pi(\cdot | s_t)$. (Policy function.)
- States are randomly sampled: $s_{t+1} \sim p(\cdot | s_t, a_t)$. (State transition.)

State-Value Function

Definition: Discounted return (aka cumulative discounted future reward).

- $R_t = r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \gamma^3 r_{t+3} + \dots$ (to infinity.)

Definition: Action-value function.

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

Definition: State-value function.

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



- Integrate out action a .
- Given s_t and π , state-value function can tell the expected return.

State-Value Function

Definition: State-value function.

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

Policy Function Approximation

Definition: State-value function.

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

Policy-based learning: Learn a policy π that maximizes $\mathbb{E}_s [V_{\pi}(s)]$.

Policy Function Approximation

Definition: State-value function.

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

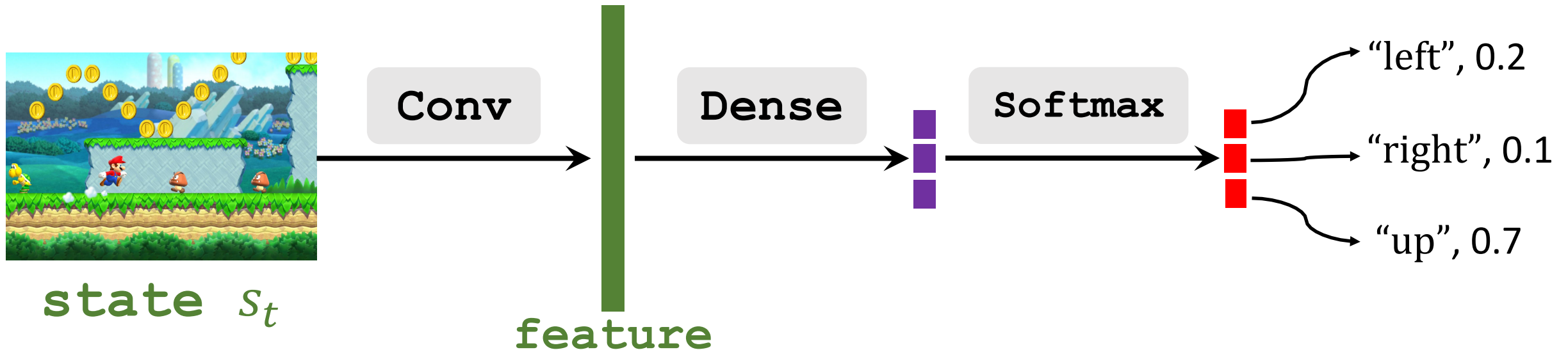
Policy-based learning: Learn a policy π that maximizes $\mathbb{E}_s [V_{\pi}(s)]$.

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

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

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

- $\pi(a|s; \theta) = 0.2$ means that observing s , the agent shall take action a with probability 0.2.
- 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.)



Policy Gradient

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 .

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



Push the differentiation into 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 .

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

Q_π is independent of θ .



Policy Gradient

Definition: Approximate state-value function.

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

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

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

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

- Thus $\frac{\partial f(x)}{\partial x} = f(x) \cdot \frac{\partial \log[f(x)]}{\partial x}.$

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. $\boldsymbol{\theta}$.

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

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

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. $\boldsymbol{\theta}$.

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

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 .

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

- Policy gradient **ascent**:

$$\theta_{t+1} \leftarrow \theta_t + \beta \cdot \frac{\partial V(s_t; \theta)}{\partial \theta} \Big|_{\theta = \theta_t}.$$

- **Increasing** the state-value means improving the policy.

Policy Gradient

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

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

Question: How to compute the policy gradient $\frac{\partial V(s; \theta)}{\partial \theta}$?

Policy Gradient

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

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

Question: How to compute the policy gradient $\frac{\partial V(s; \theta)}{\partial \theta}$?

- Sample actions: $a^{(1)}, a^{(2)}, \dots, a^{(k)} \sim \pi(\cdot | s; \theta)$.
 - (The agent does not actually perform the actions.)
 - Sample only one action (i.e., $k = 1$) also works.
- Compute $\tilde{\mathbf{g}}(\theta) = \frac{1}{k} \sum_{i=1}^k \frac{\partial \log \pi(a^{(i)} | s_t, \theta)}{\partial \theta} \cdot Q_{\pi}(s_t, a^{(i)})$.
- $\tilde{\mathbf{g}}(\theta)$ is unbiased estimate of $\frac{\partial V(s; \theta)}{\partial \theta}$.

Update policy network using policy gradient

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



The agent may not actually perform action a_t .

Update policy network using policy gradient

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. (Stochastic) policy gradient: $\tilde{\mathbf{g}}(\theta_t) \approx q_t \cdot \mathbf{d}_{\theta,t}$.
6. Update policy network: $\theta_{t+1} = \theta_t + \beta \cdot \tilde{\mathbf{g}}(\theta_t)$.

Update policy network using policy gradient

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. (Stochastic) policy gradient: $\tilde{\mathbf{g}}(\theta_t) \approx q_t \cdot \mathbf{d}_{\theta,t}$.
6. Update policy network: $\theta_{t+1} = \theta_t + \beta \cdot \tilde{\mathbf{g}}(\theta_t)$.

Update policy network using policy gradient

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: Monte Carlo.

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

$$s_t, a_t, r_t, s_{t+1}, a_{t+1}, r_{t+1}, \dots, s_T, a_T, r_T.$$

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

Update policy network using policy gradient

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.

Summary

Policy-Based Method

- If a good policy function $\pi(a|s)$ is known, the agent can follow the policy: $a_t \sim \pi(\cdot | s_t)$.
- Approximate policy function $\pi(a|s)$ by policy network $\pi(a|s; \theta)$.
- Learn the policy network by policy gradient.
- Policy gradient algorithm learn θ that maximizes $\mathbb{E}_s[V(s; \theta)]$.