

REINFORCE with Baseline

Shusen Wang

Policy Gradient with Baseline

Value Functions

- Discounted return:

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

- Action-value function:

$$Q_\pi(s_t, a_t) = \mathbb{E}[U_t \mid s_t, a_t].$$

- State-value function:

$$V_\pi(s_t) = \mathbb{E}_{\mathbf{A}}[Q_\pi(s_t, \mathbf{A}) \mid s_t].$$

Policy Gradient with Baseline

- Use policy network, $\pi(\textcolor{red}{a}|\textcolor{green}{s}; \boldsymbol{\theta})$, for controlling the agent.
- State-value function:

$$\begin{aligned}\underline{V_{\pi}(\textcolor{green}{s})} &= \mathbb{E}_{\textcolor{red}{A} \sim \pi}[Q_{\pi}(\textcolor{green}{s}, \textcolor{red}{A})] \\ &= \sum_{\textcolor{red}{a}} \underline{\pi(\textcolor{red}{a}|\textcolor{green}{s}; \boldsymbol{\theta})} \cdot Q_{\pi}(\textcolor{green}{s}, \textcolor{red}{a}).\end{aligned}$$

- Policy gradient with baseline:

$$\frac{\partial V_{\pi}(\textcolor{green}{s})}{\partial \boldsymbol{\theta}} = \mathbb{E}_{\textcolor{red}{A} \sim \pi} \left[\underline{\frac{\partial \ln \pi(\textcolor{red}{A} | \textcolor{green}{s}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}}} \cdot \underline{(Q_{\pi}(\textcolor{green}{s}, \textcolor{red}{A}) - V_{\pi}(\textcolor{green}{s}_t))} \right].$$

Policy Gradient with Baseline

Policy gradient with baseline:

$$\frac{\partial V_{\pi}(s_t)}{\partial \theta} = \mathbb{E}_{A_t \sim \pi} \left[\frac{\partial \ln \pi(A_t | s_t; \theta)}{\partial \theta} \cdot (Q_{\pi}(s_t, A_t) - V_{\pi}(s_t)) \right].$$

Policy Gradient with Baseline

Policy gradient with baseline:

$$\frac{\partial V_{\pi}(s_t)}{\partial \theta} = \mathbb{E}_{A_t \sim \pi} \left[\frac{\partial \ln \pi(A_t | s_t; \theta)}{\partial \theta} \cdot (Q_{\pi}(s_t, A_t) - V_{\pi}(s_t)) \right].$$
$$= \mathbf{g}(A_t)$$

Policy Gradient with Baseline

Policy gradient with baseline:

$$\frac{\partial V_{\pi}(s_t)}{\partial \theta} = \mathbb{E}_{A_t \sim \pi} \left[\frac{\partial \ln \pi(A_t | s_t; \theta)}{\partial \theta} \cdot (Q_{\pi}(s_t, A_t) - V_{\pi}(s_t)) \right].$$
$$= \mathbf{g}(A_t)$$

- Randomly sample $a_t \sim \pi(\cdot | s_t; \theta)$.
- Then $\mathbf{g}(a_t)$ is an unbiased estimation of the policy gradient.

Approximations

Policy gradient with baseline:

$$\frac{\partial V_{\pi}(s_t)}{\partial \theta} = \mathbb{E}_{A_t \sim \pi} \left[\frac{\partial \ln \pi(A_t | s_t; \theta)}{\partial \theta} \cdot (Q_{\pi}(s_t, A_t) - V_{\pi}(s_t)) \right].$$
$$= \mathbf{g}(A_t)$$

Stochastic policy gradient with baseline:

$$\mathbf{g}(a_t) = \frac{\partial \ln \pi(a_t | s_t; \theta)}{\partial \theta} \cdot (Q_{\pi}(s_t, a_t) - V_{\pi}(s_t)).$$

Approximations

Stochastic policy gradient with baseline:

$$\mathbf{g}(a_t) = \frac{\partial \ln \pi(a_t | s_t; \theta)}{\partial \theta} \cdot (Q_\pi(s_t, a_t) - V_\pi(s_t)).$$

Approximations

Stochastic policy gradient with baseline:

$$\mathbf{g}(a_t) = \frac{\partial \ln \pi(a_t | s_t; \theta)}{\partial \theta} \cdot (Q_{\pi}(s_t, a_t) - V_{\pi}(s_t))$$

Approximations

Stochastic policy gradient with baseline:

$$\mathbf{g}(a_t) = \frac{\partial \ln \pi(a_t | s_t; \theta)}{\partial \theta} \cdot (Q_\pi(s_t, a_t) - V_\pi(s_t)).$$

- Recall that $Q_\pi(s_t, a_t) = \mathbb{E}[U_t \mid s_t, a_t]$.
- Monte Carlo approximation to $Q_\pi(s_t, a_t) \approx u_t$ (REINFORCE):
 - Observing 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 return: $u_t = \sum_{i=t}^T \gamma^{i-t} \cdot r_i$.
 - u_t is unbiased Monte Carlo estimate of $Q_\pi(s_t, a_t)$.

Approximations

Stochastic policy gradient with baseline:

$$\mathbf{g}(a_t) = \frac{\partial \ln \pi(a_t | s_t; \theta)}{\partial \theta} \cdot (Q_\pi(s_t, a_t) - V_\pi(s_t))$$

- Approximate $V(s; \theta)$ by the value network, $v(s; \mathbf{w})$.

Approximations

Approximate policy gradient with baseline:

$$\frac{\partial V_{\pi}(s_t)}{\partial \theta} \approx \mathbf{g}(a_t) \approx \frac{\partial \ln \pi(a_t | s_t; \theta)}{\partial \theta} \cdot (u_t - v(s_t; \mathbf{w})).$$

- Three approximations:
 1. Approximate expectation using one sample, a_t . (Monte Carlo.)
 2. Approximate $Q_{\pi}(s_t, a_t)$ by u_t . (Another Monte Carlo.)
 3. Approximate $V_{\pi}(s)$ by the value network, $v(s; \mathbf{w})$.

Summary of Approximations

$$\frac{\partial V_{\pi}(s_t)}{\partial \theta} = \mathbb{E}_{A_t \sim \pi} \left[\frac{\partial \ln \pi(A_t | s_t; \theta)}{\partial \theta} \cdot (Q_{\pi}(s_t, A_t) - V_{\pi}(s_t)) \right].$$



$$\mathbf{g}(a_t) = \frac{\partial \ln \pi(a_t | s_t; \theta)}{\partial \theta} \cdot (Q_{\pi}(s_t, a_t) - V_{\pi}(s_t)).$$

Summary of Approximations

$$\frac{\partial V_{\pi}(s_t)}{\partial \theta} = \mathbb{E}_{A_t \sim \pi} \left[\frac{\partial \ln \pi(A_t | s_t; \theta)}{\partial \theta} \cdot (Q_{\pi}(s_t, A_t) - V_{\pi}(s_t)) \right].$$



$$\mathbf{g}(a_t) = \frac{\partial \ln \pi(a_t | s_t; \theta)}{\partial \theta} \cdot (Q_{\pi}(s_t, a_t) - V_{\pi}(s_t)).$$

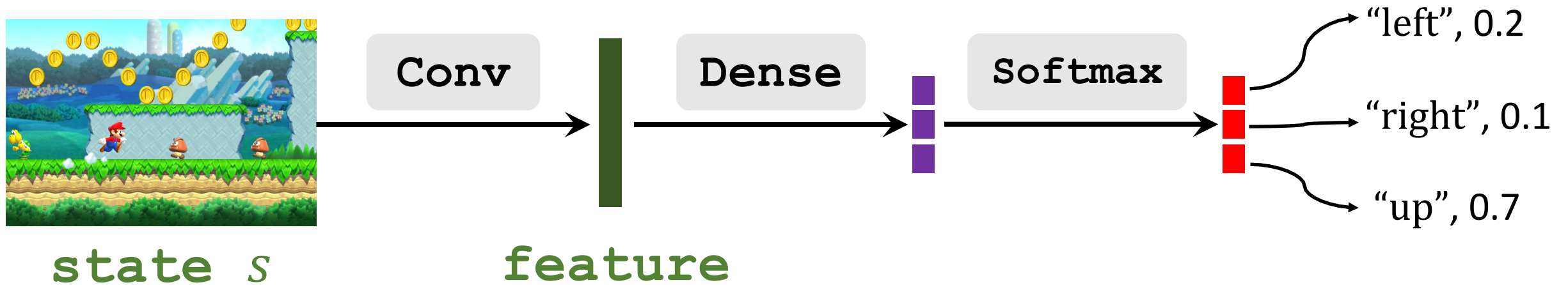


$$\mathbf{g}(a_t) \approx \frac{\partial \ln \pi(a_t | s_t; \theta)}{\partial \theta} \cdot (u_t - v(s_t; \mathbf{w})).$$

Policy and Value Networks

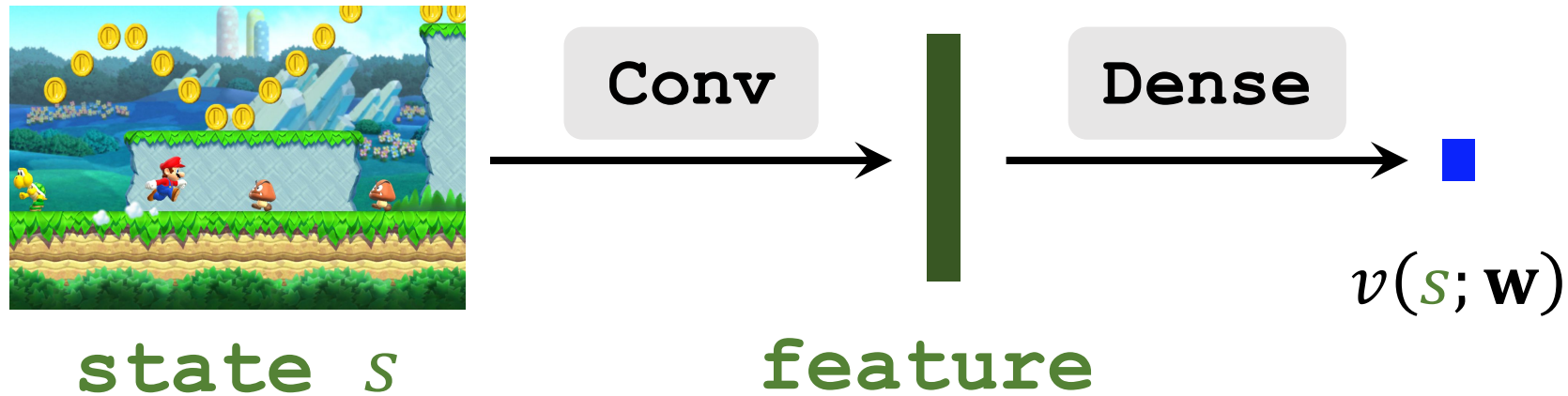
Policy Network

Approximate policy function, $\pi(a|s)$, by policy network, $\pi(a|s; \theta)$.

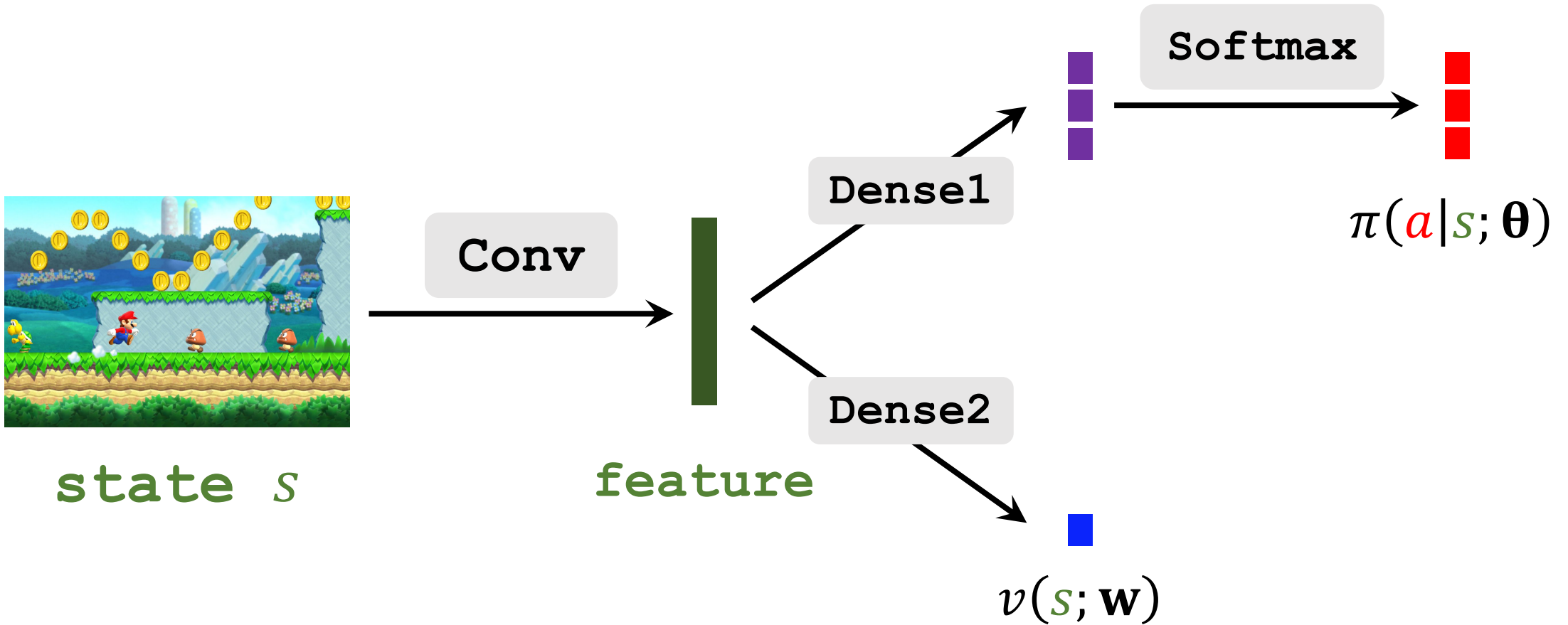


Value Network

Approximate state-value, $V_{\pi}(s)$, by value network, $v(s; \mathbf{w})$.



Parameter Sharing



REINFORCE with Baseline

Updating the policy network

Approximate policy gradient with baseline:

$$\frac{\partial V_{\pi}(s_t)}{\partial \theta} \approx \frac{\partial \ln \pi(a_t | s_t; \theta)}{\partial \theta} \cdot (u_t - v(s_t; \mathbf{w})).$$

- Update policy network by policy gradient ascent:

$$\theta \leftarrow \theta + \beta \cdot \frac{\partial \ln \pi(a_t | s_t; \theta)}{\partial \theta} \cdot (u_t - v(s_t; \mathbf{w})).$$

Updating the value network

- Recall $v(s_t; \mathbf{w})$ is an approximation to $V_\pi(s_t) = \mathbb{E}[U_t | s_t]$.

Updating the value network

- Recall $v(s_t; \mathbf{w})$ is an approximation to $V(s_t; \boldsymbol{\theta}) = \mathbb{E}[U_t \mid s_t]$.
- Encourage $v(s_t; \mathbf{w})$ to approach u_t by decreasing:

$$\underline{\delta_t} = \underline{u_t} - \underline{v(s_t; \mathbf{w})}.$$

Updating the value network

- Recall $v(s_t; \mathbf{w})$ is an approximation to $V(s_t; \boldsymbol{\theta}) = \mathbb{E}[U_t \mid s_t]$.
- Encourage $v(s_t; \mathbf{w})$ to approach u_t by decreasing:

$$\delta_t = u_t - v(s_t; \mathbf{w}).$$

- Gradient: $\frac{\partial \delta_t^2/2}{\partial \mathbf{w}} = -\delta_t \cdot \frac{\partial v(s_t; \mathbf{w})}{\partial \mathbf{w}}.$

- Gradient descent:

$$\mathbf{w} \leftarrow \mathbf{w} - \alpha \cdot (-\delta_t) \cdot \frac{\partial v(s_t; \mathbf{w})}{\partial \mathbf{w}}.$$

Summary of Algorithm

- Play a game to the end and observe the trajectory:

$$s_1, a_1, r_1, s_2, a_2, r_2, \dots, s_n, a_n, r_n .$$

- Compute $u_t = \sum_{i=t}^T \gamma^{i-t} \cdot r_i$ and $\delta_t = u_t - v(s_t; \mathbf{w})$.

Summary of Algorithm

- Play a game to the end and observe the trajectory:

$$s_1, a_1, r_1, s_2, a_2, r_2, \dots, s_n, a_n, r_n.$$

- Compute $u_t = \sum_{i=t}^T \gamma^{i-t} \cdot r_i$ and $\delta_t = u_t - v(s_t; \mathbf{w})$.
- Update the policy network by:

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \beta \cdot \delta_t \cdot \frac{\partial \ln \pi(a_t | s_t; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}}.$$

- Update the value network by:

$$\mathbf{w} \leftarrow \mathbf{w} + \alpha \cdot \delta_t \cdot \frac{\partial v(s_t; \mathbf{w})}{\partial \mathbf{w}}.$$

Summary of Algorithm

- Play a game to the end and observe the trajectory:

$$s_1, a_1, r_1, s_2, a_2, r_2, \dots, s_n, a_n, r_n.$$

- Compute $u_t = \sum_{i=t}^T \gamma^{i-t} \cdot r_i$ and $\delta_t = u_t - v(s_t; \mathbf{w})$.
- Update the policy network by:

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \beta \cdot \delta_t \cdot \frac{\partial \ln \pi(a_t | s_t; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}}.$$

- Update the value network by:

$$\mathbf{w} \leftarrow \mathbf{w} + \alpha \cdot \delta_t \cdot \frac{\partial v(s_t; \mathbf{w})}{\partial \mathbf{w}}.$$

Summary of Algorithm

- Play a game to the end and observe the trajectory:

$$s_1, a_1, r_1, s_2, a_2, r_2, \dots, s_n, a_n, r_n.$$

- Compute $u_t = \sum_{i=t}^T \gamma^{i-t} \cdot r_i$ and $\delta_t = u_t - v(s_t; \mathbf{w})$.
- Update the policy network by:

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \beta \cdot \delta_t \cdot \frac{\partial \ln \pi(a_t | s_t; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}}.$$

- Update the value network by:

$$\mathbf{w} \leftarrow \mathbf{w} + \alpha \cdot \delta_t \cdot \frac{\partial v(s_t; \mathbf{w})}{\partial \mathbf{w}}.$$

Repeat this procedure for all t from 1 to n .

Thank you!