

Dueling Network

Shusen Wang

<http://wangshusen.github.io/>

Advantage Function

Return

Definition: Discounted return.

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

Value Functions

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.

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

Value Functions

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.

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

Optimal Value Functions

Definition: Optimal action-value function.

- $Q^*(s, a) = \max_{\pi} Q_{\pi}(s, a).$

Optimal Value Functions

Definition: Optimal action-value function.

- $Q^*(s, a) = \max_{\pi} Q_{\pi}(s, a).$

Definition: Optimal state-value function.

- $V^*(s) = \max_{\pi} V_{\pi}(s).$

Optimal Value Functions

Definition: Optimal action-value function.

- $Q^*(s, a) = \max_{\pi} Q_{\pi}(s, a).$

Definition: Optimal state-value function.

- $V^*(s) = \max_{\pi} V_{\pi}(s).$

Definition: Optimal advantage function.

- $A^*(s, a) = \underline{Q^*(s, a)} - \underline{V^*(s)}.$

Properties of Advantage Function

Theorem 1: $V^*(s) = \max_a Q^*(s, a)$.

- Recall the definition of the optimal advantage function:

$$A^*(s, a) = Q^*(s, a) - V^*(s).$$

- It follows that

$$\max_a A^*(s, a) = \max_a Q^*(s, a) - V^*(s).$$

Properties of Advantage Function

Theorem 1: $V^*(s) = \max_a Q^*(s, a)$.

- Recall the definition of the optimal advantage function:

$$A^*(s, a) = Q^*(s, a) - V^*(s).$$

- It follows that

$$\max_a A^*(s, a) = \max_a Q^*(s, a) - V^*(s).$$

- Using Theorem 1, we get $\max_a A^*(s, a) = 0$.

Properties of Advantage Function

Definition of advantage: $A^*(s, a) = Q^*(s, a) - V^*(s)$.



Theorem 2: $Q^*(s, a) = V^*(s) + A^*(s, a)$

Properties of Advantage Function

Definition of advantage: $A^*(s, a) = Q^*(s, a) - V^*(s)$.



Theorem 2: $Q^*(s, a) = V^*(s) + A^*(s, a) - \max_a A^*(s, a).$
 $= 0$

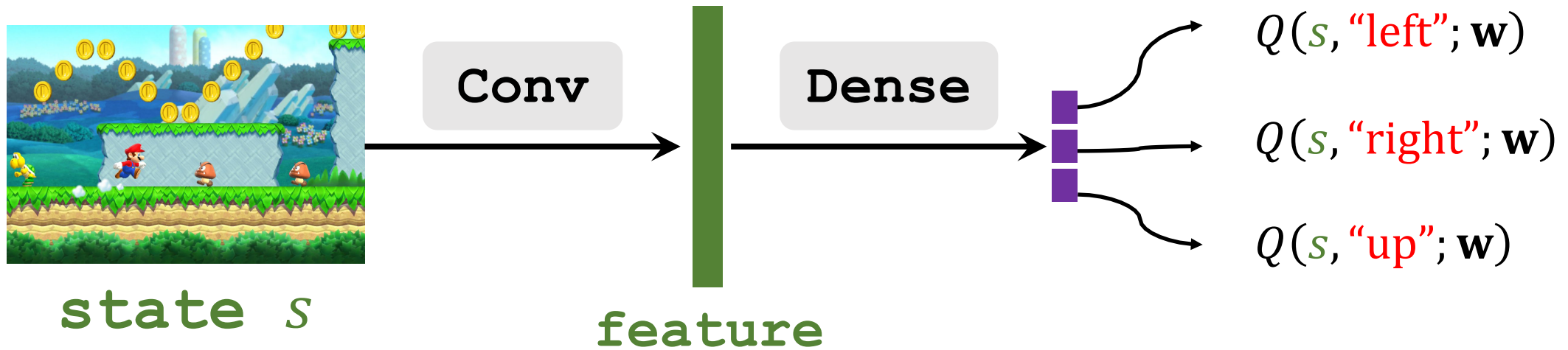
Dueling Network

Reference:

1. Wang et al. [Dueling network architectures for deep reinforcement learning](#). In *ICML*, 2016.

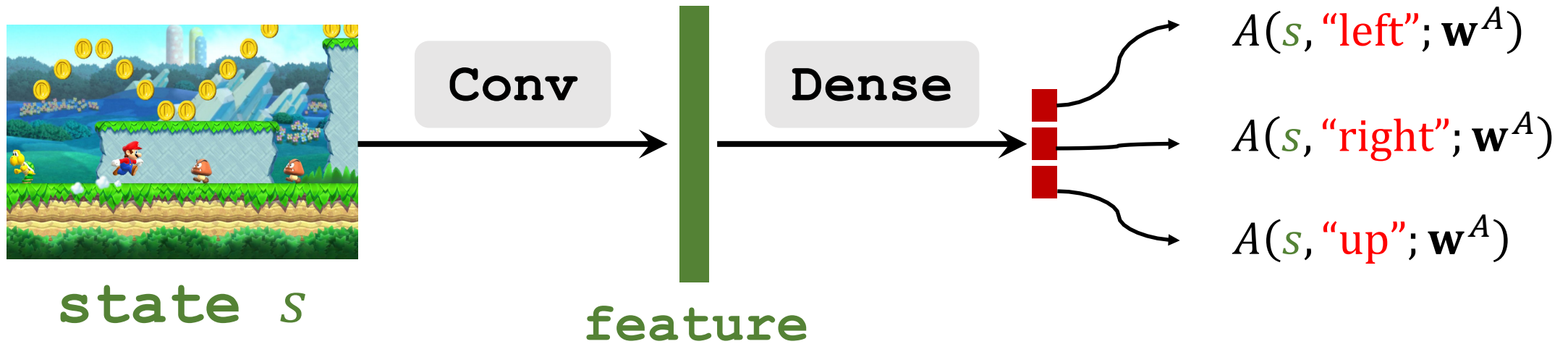
Revisiting DQN

- Approximate $Q^*(s, a)$ by a neural network, $Q(s, a; \mathbf{w})$.



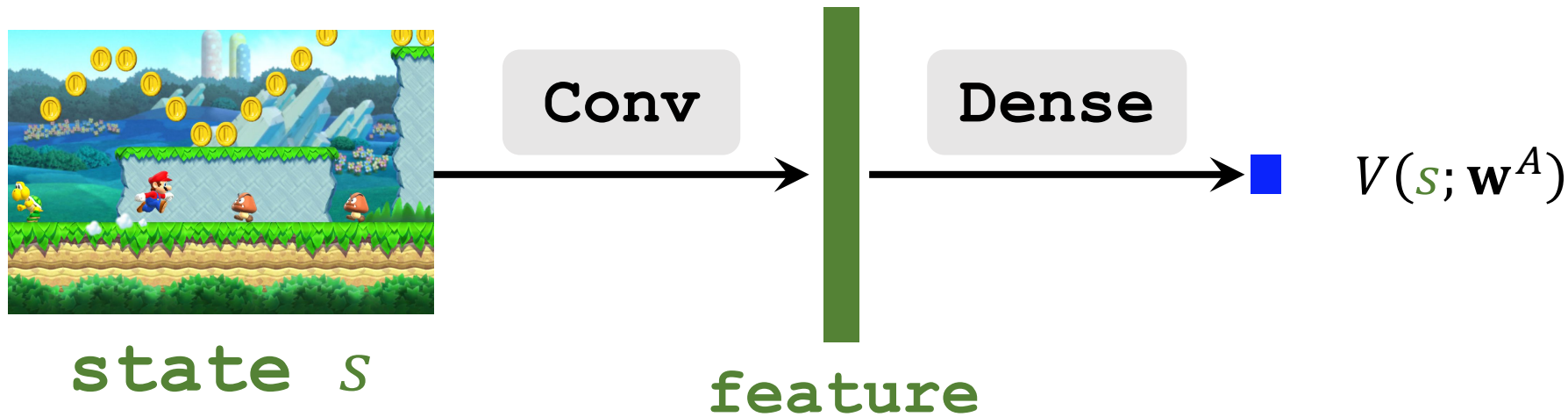
Approximating Advantage Function

- Approximate $A^*(s, a)$ by a neural network, $A(s, a; \mathbf{w}^A)$.



Approximating State-Value Function

- Approximate $V^*(s)$ by a neural network, $V(s; \mathbf{w}^V)$.



Dueling Network

Theorem 2: $Q^*(s, a) = V^*(s) + A^*(s, a) - \max_a A^*(s, a).$

- Approximate $V^*(s)$ by a neural network, $V(s; \mathbf{w}^V).$

Dueling Network

Theorem 2: $Q^*(s, a) = V^*(s) + A^*(s, a) - \max_a A^*(s, a).$

- Approximate $V^*(s)$ by a neural network, $V(s; \mathbf{w}^V)$.
- Approximate $A^*(s, a)$ by a neural network, $A(s, a; \mathbf{w}^A)$.

Dueling Network

Theorem 2: $Q^*(s, a) = V^*(s) + A^*(s, a) - \max_a A^*(s, a)$.

- Approximate $V^*(s)$ by a neural network, $V(s; \mathbf{w}^V)$.
- Approximate $A^*(s, a)$ by a neural network, $A(s, a; \mathbf{w}^A)$.
- Thus, approximate $Q^*(s, a)$ by the dueling network:

$$Q(s, a; \mathbf{w}^A, \mathbf{w}^V) = V(s; \mathbf{w}^V) + A(s, a; \mathbf{w}^A) - \max_a A(s, a; \mathbf{w}^A).$$

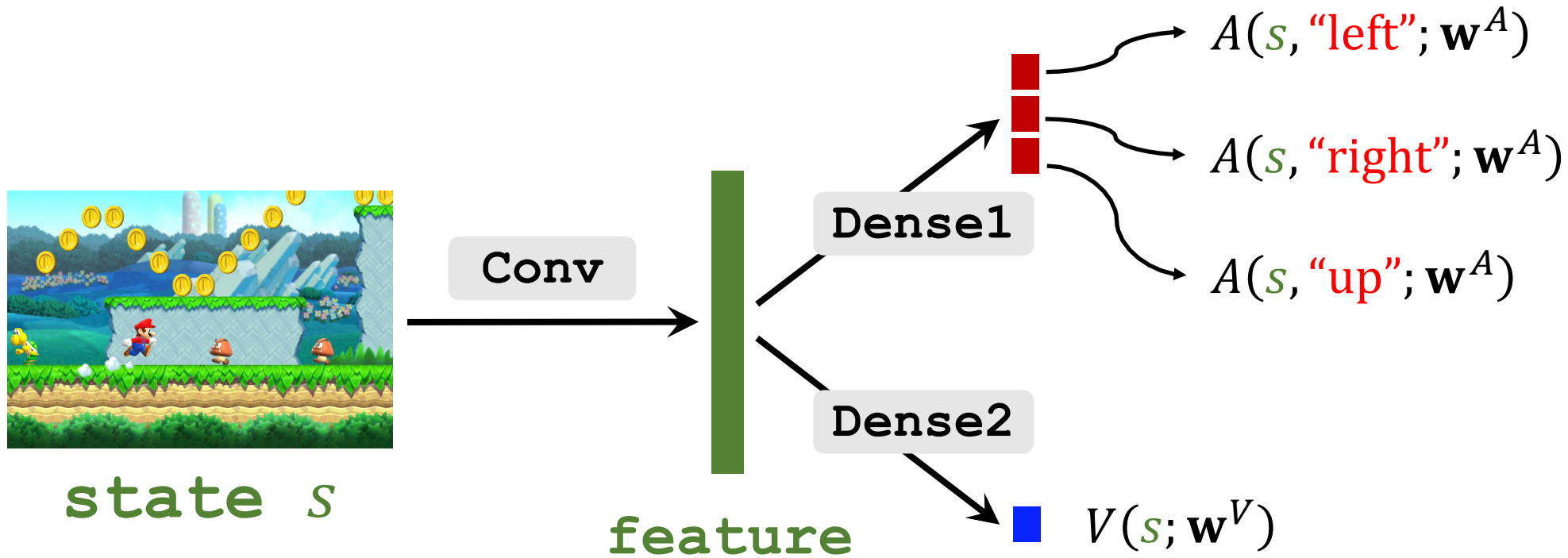
Dueling Network

$$Q(s, a; \mathbf{w}) = V(s; \mathbf{w}^V) + A(s, a; \mathbf{w}^A) - \max_a A(s, a; \mathbf{w}^A).$$

Here, $\mathbf{w} = (\mathbf{w}^A, \mathbf{w}^V)$

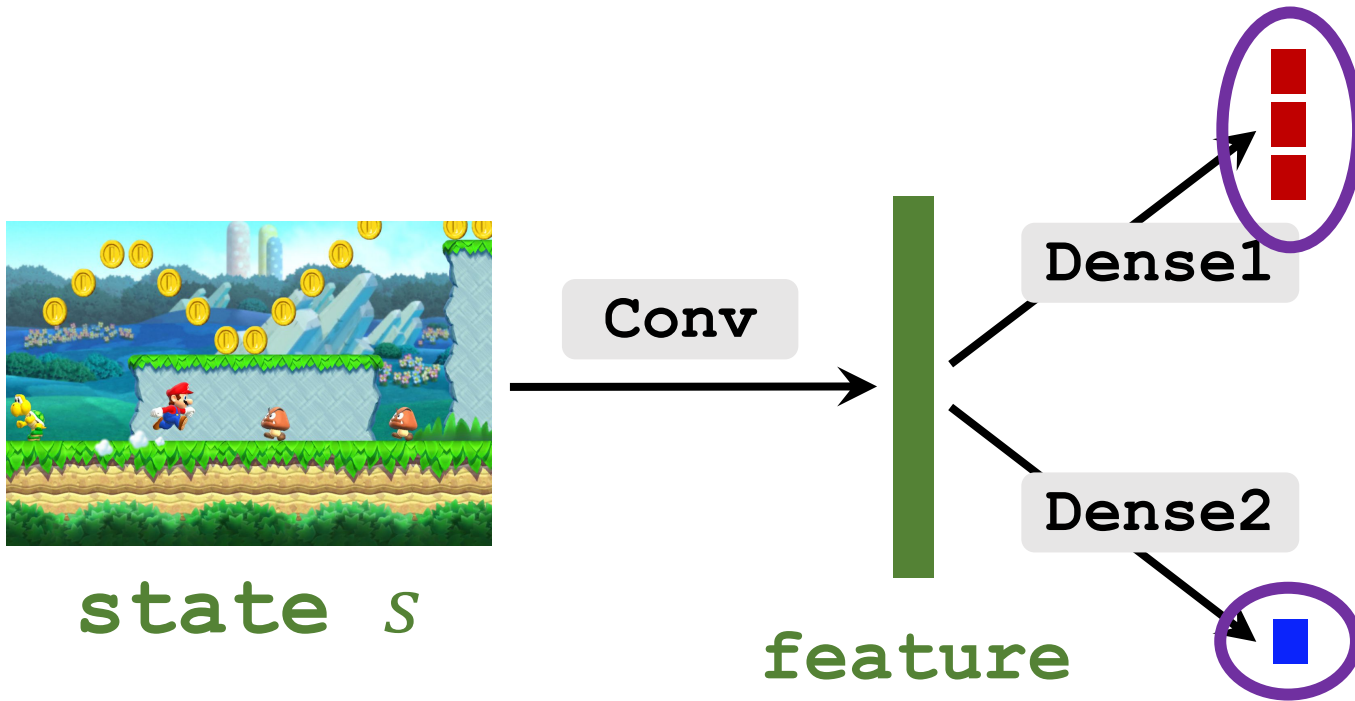
Dueling Network

$$Q(s, a; \mathbf{w}) = V(s; \mathbf{w}^V) + A(s, a; \mathbf{w}^A) - \max_a A(s, a; \mathbf{w}^A).$$



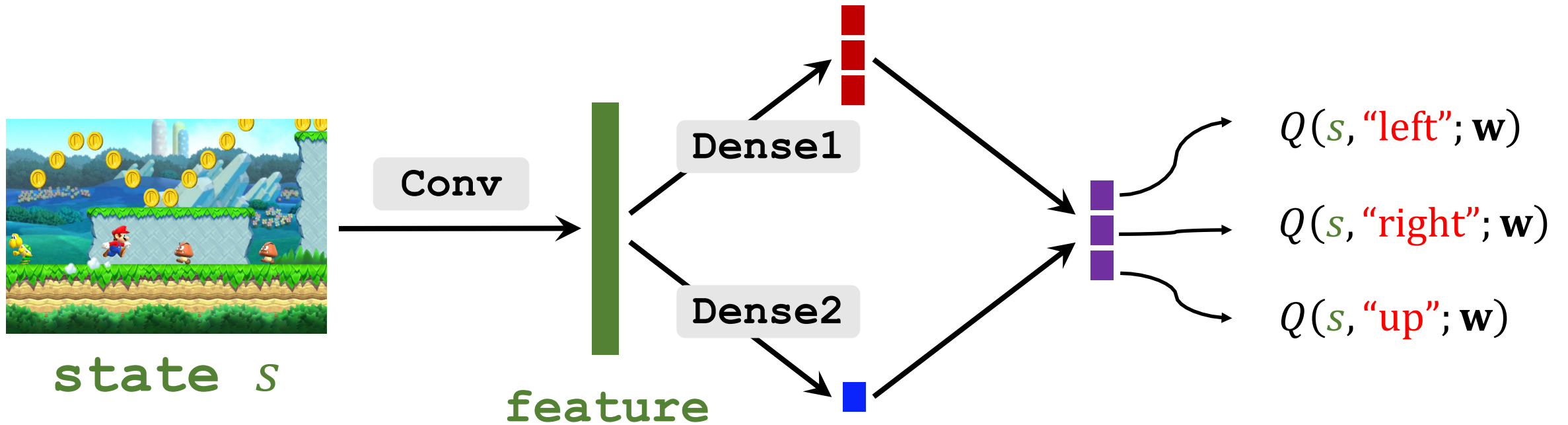
Dueling Network

$$Q(s, a; \mathbf{w}) = V(s; \mathbf{w}^V) + A(s, a; \mathbf{w}^A) - \max_a A(s, a; \mathbf{w}^A).$$



Dueling Network

$$Q(s, a; \mathbf{w}) = V(s; \mathbf{w}^V) + A(s, a; \mathbf{w}^A) - \max_a A(s, a; \mathbf{w}^A).$$



Training

- Dueling network, $Q(s, a; \mathbf{w})$, is an approximation to $Q^*(s, a)$.
- Learn the parameter, $\mathbf{w} = (\mathbf{w}^A, \mathbf{w}^V)$, in the same way as the other DQNs.
- Tricks can be used in the same way.
 - Prioritized experience replay.
 - Double DQN.
 - Multi-step TD target.

Overcome Non-identifiability

Problem of Non-identifiability

➡ • **Equation 1:** $Q^*(s, a) = V^*(s) + A^*(s, a).$

➡ • **Equation 2:** $Q^*(s, a) = V^*(s) + A^*(s, a) - \max_a A^*(s, a)$

Question: Why is the zero term necessary?

Problem of Non-identifiability

- **Equation 1:** $Q^*(s, a) = V^*(s) + A^*(s, a)$.
- Equation 1 has the problem of *non-identifiability*.
 - Let $V' = V^* + 10$ and $A' = A^* - 10$.
 - Then $Q^*(s, a) = V^*(s) + A^*(s, a) = V'(s) + A'(s, a)$.
- Why is non-identifiability a problem?

Problem of Non-identifiability

- **Equation 2:** $Q^*(s, a) = V^*(s) + A^*(s, a) - \max_a A^*(s, a)$.
- Equation 2 does not have the problem.

Dueling Network

$$Q(s, a; \mathbf{w}) = V(s; \mathbf{w}^V) + A(s, a; \mathbf{w}^A) - \max_a A(s, a; \mathbf{w}^A).$$

Alternative:

$$Q(s, a; \mathbf{w}) = V(s; \mathbf{w}^V) + A(s, a; \mathbf{w}^A) - \text{mean}_a A(s, a; \mathbf{w}^A).$$

Thank you!