

State-Value Function Approximation

Definition: State-value function.

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

Policy network (actor):

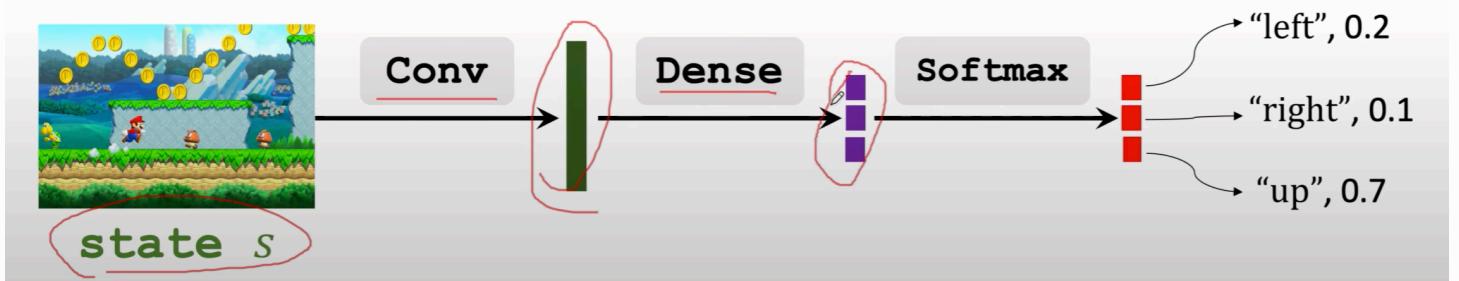
- Use neural net $\pi(a|s; \theta)$ to approximate $\pi(a|s)$.
- θ : trainable parameters of the neural net.

Value network (critic):

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

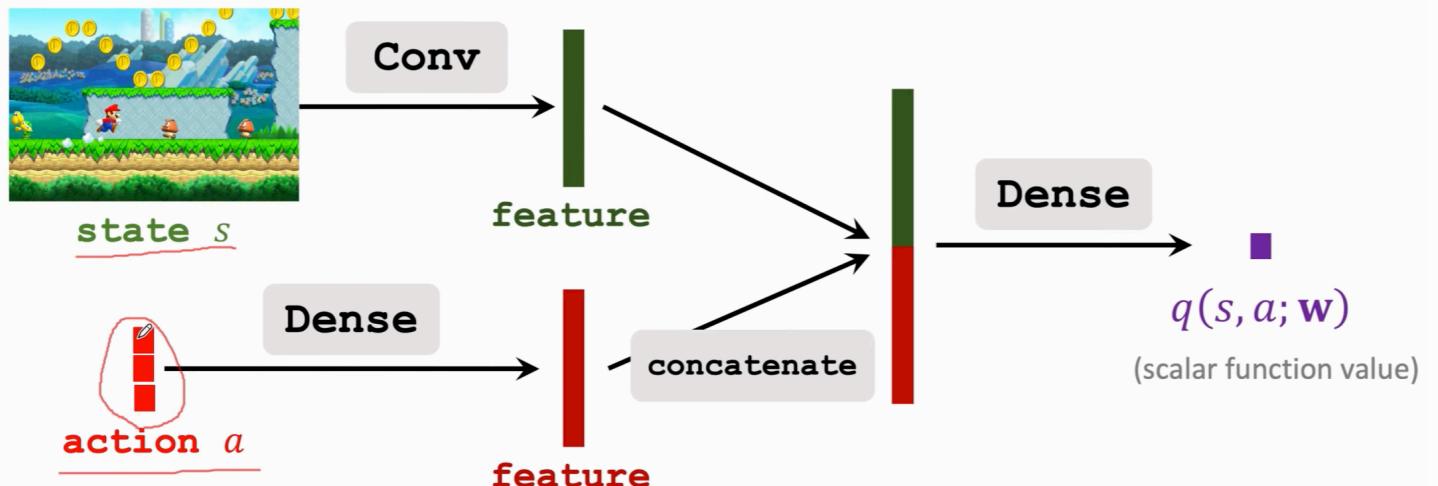
Policy Network (Actor): $\pi(a|s, \theta)$

- 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): $q(s, a; w)$

- Inputs: state s and action a .
- Output: approximate action-value (scalar).



Train the networks

Definition: State-value function approximated using neural networks.

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

Training: Update the parameters θ and w .

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

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

Training: Update the parameters θ and w .

1. Observe the state s_t .
2. Randomly sample action a_t according to $\pi(\cdot | s_t; \theta_t)$.
3. Perform a_t and observe new state s_{t+1} and reward r_t .
4. Update w (in value network) using temporal difference (TD).
5. Update θ (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: $y_t = r_t + \gamma \cdot q(s_{t+1}, a_{t+1}; \mathbf{w}_t)$.
-
- 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.

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

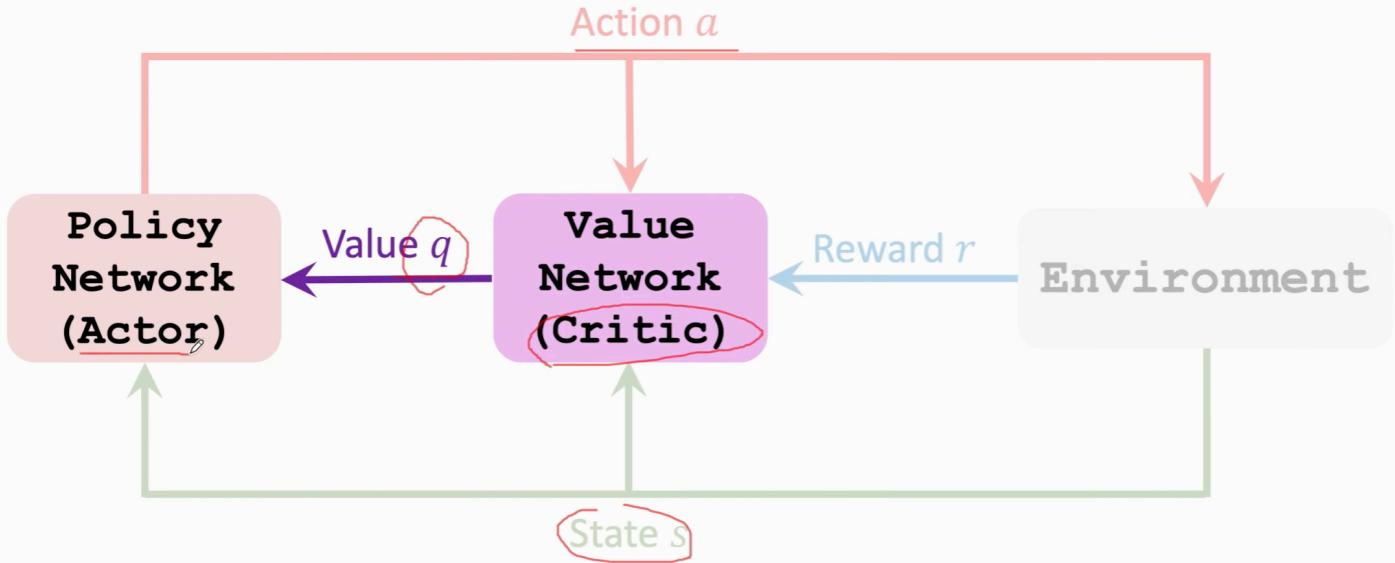
Policy gradient: Derivative of $V(s_t; \theta, \mathbf{w})$ w.r.t. θ .

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

Algorithm: Update policy network using stochastic policy gradient.

- Random sampling: $a \sim \pi(\cdot | s_t; \theta_t)$. (Thus $\mathbf{g}(a, \theta)$ is unbiased.)
- Stochastic gradient ascent: $\theta_{t+1} = \theta_t + \beta \cdot \mathbf{g}(a, \theta_t)$.

Actor-Critic Method



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: $\mathbf{d}_{w,t} = \frac{\partial q(s_t, a_t; \mathbf{w})}{\partial \mathbf{w}} |_{\mathbf{w}=\mathbf{w}_t}$.
7. Update value network: $\mathbf{w}_{t+1} = \mathbf{w}_t - \alpha \cdot \delta_t \cdot \mathbf{d}_{w,t}$.
8. Differentiate policy network: $\mathbf{d}_{\theta,t} = \frac{\partial \log \pi(a_t | s_t, \theta)}{\partial \theta} |_{\theta=\theta_t}$.
9. Update policy network: $\theta_{t+1} = \theta_t + \beta \cdot q_t \cdot \mathbf{d}_{\theta,t}$.

Policy Gradient with Baseline

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: $\mathbf{d}_{w,t} = \frac{\partial q(s_t, a_t; \mathbf{w})}{\partial \mathbf{w}} |_{\mathbf{w}=\mathbf{w}_t}$.
7. Update value network: $\mathbf{w}_{t+1} = \mathbf{w}_t - \alpha \cdot \delta_t \cdot \mathbf{d}_{w,t}$.
8. Differentiate policy network: $\mathbf{d}_{\theta,t} = \frac{\partial \log \pi(a_t | s_t, \theta)}{\partial \theta} |_{\theta=\theta_t}$.
9. Update policy network: $\theta_{t+1} = \theta_t + \beta \cdot \delta_t \cdot \mathbf{d}_{\theta,t}$.

Policy Network and Value Network

Definition: State-value function.

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

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.

After training

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

Training

Learning: Update the policy network (actor) by policy gradient.

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

Learning: Update the value network (critic) by TD learning.

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

