



澳門大學
UNIVERSIDADE DE MACAU
UNIVERSITY OF MACAU

Reinforcement Learning

GongYing

gong.ying@connect.umac.mo

2023.12.08

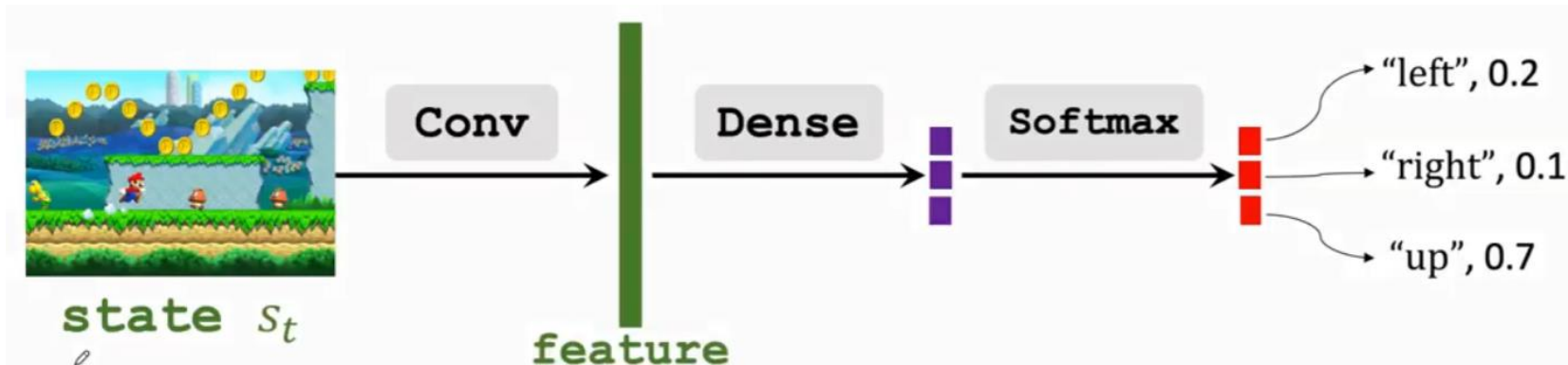
Contents

- **Policy-based reinforcement learning**
 - Policy gradient
 - REINFORCE & actor-critic
- **Policy gradient with baseline**
 - REINFORCE with baseline
 - A2C: advantage actor-critic
- **Advanced techniques in policy-based RL**
 - TRPO: trust region policy optimization
 - Entropy regularization
- **Partial observations**
 - RNN: recurrent neural network



Policy gradient

- **Policy function:** $\pi(a|s)$, $\sum_{a \in A} \pi(a|s) = 1$.
- **Input:** s .
- **Output:** probabilities for all actions, e.g.,
 $\pi(\text{left}|s)=0.2$,
 $\pi(\text{right}|s)=0.1$,
 $\pi(\text{up}|s)=0.7$.
- **Policy network:** use $\pi(a|s;\theta)$ to approximate $\pi(a|s)$.



Policy gradient

- **Discounted return:**

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

- **Action-value function:**

$$Q_\pi(s_t, a_t) = E[U_t | S_t = s_t, A_t = a_t]$$

- **State-value function:**

$$V_\pi(s_t) = E_A[Q_\pi(s_t, A)]$$

$$\text{(If discrete)} = \sum_a \pi(a | s_t) \cdot Q_\pi(s_t, a).$$

- **Approximate state-value function:**

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

Policy gradient

- **Approximate state-value function:**

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

If a policy is good, the mean value of $V_\pi(s_t)$ is supposed to be big, so we can define the target function:

$$J(\theta) = E_S[V_\pi(S)].$$

The target is to maximize $J(\theta)$ by updating θ :

$$\theta \leftarrow \theta + \beta \cdot \frac{\partial V(s; \theta)}{\partial \theta}$$

policy gradient

Policy gradient

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

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

(assume that $Q_{\pi}(s, a)$ is independent of θ)

$$= \sum_a \frac{\partial \pi(a|s; \theta)}{\partial \theta} \cdot Q_{\pi}(s, a) \quad \text{--- Form 1}$$

$$= \sum_a \left\{ \pi(a|s; \theta) \cdot \left[\frac{\partial \ln \pi(a|s; \theta)}{\partial \theta} \cdot Q_{\pi}(s, a) \right] \right\}$$

$$= \mathbb{E}_A \left[\frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} \cdot Q_{\pi}(s, A) \right] \quad \text{--- Form 2}$$



Policy gradient

2 forms of policy gradient:

- Form 1: $\frac{\partial V(s; \theta)}{\partial \theta} = \sum_a \frac{\partial \pi(a|s; \theta)}{\partial \theta} \cdot Q_{\pi}(s, a)$
- Form 2: $\frac{\partial V(s; \theta)}{\partial \theta} = E_{A \sim \pi(\cdot|s; \theta)} \left[\frac{\partial \log \pi(A|s; \theta)}{\partial \theta} \cdot Q_{\pi}(s, A) \right]$



Policy gradient

- Use form 1: $\frac{\partial V(s; \theta)}{\partial \theta} = \sum_a \frac{\partial \pi(a|s; \theta)}{\partial \theta} \cdot Q_{\pi}(s, a)$
 1. For every $a \in A$, calculate $f(a, \theta) = \frac{\partial \pi(a|s; \theta)}{\partial \theta} \cdot Q_{\pi}(s, a)$.
 2. $\frac{\partial V(s; \theta)}{\partial \theta} = \sum_a f(a, \theta)$.
- E.g., $A = \{\text{"left"}, \text{"right"}, \text{"up"}\}$.

$$\frac{\partial V(s; \theta)}{\partial \theta} = f(\text{"left"}, \theta) + f(\text{"right"}, \theta) + f(\text{"up"}, \theta).$$



Policy gradient

- Use form 2: $\frac{\partial V(s; \theta)}{\partial \theta} = E_{A \sim \pi(\cdot | s; \theta)} \left[\frac{\partial \log \pi(A | s; \theta)}{\partial \theta} \cdot Q_{\pi}(s, A) \right]$
 1. (Monte Carlo) Randomly sample \hat{a} according to $\pi(\cdot | s; \theta)$.
 2. Calculate $g(\hat{a}, \theta) = \frac{\partial \log \pi(\hat{a} | s; \theta)}{\partial \theta} \cdot Q_{\pi}(s, \hat{a})$.
 3. Use $g(\hat{a}, \theta)$ to approximate $\frac{\partial V(s; \theta)}{\partial \theta}$.



Policy gradient

Algorithm:

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



REINFORCE & actor-critic

REINFORCE: Monte Carlo

- Play **the whole game** and record 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 .
- Because $Q_\pi(s_t, a_t) = E[U_t | S_t = s_t, A_t = a_t]$, we can use u_t to approximate $Q_\pi(s_t, a_t)$.
- Use $q_t = u_t$.

REINFORCE & actor-critic

Actor-critic method: use neural network to approximate Q_π

- **State-value function**

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

- **Policy network – actor**

Use neural network $\pi(a|s; \theta)$ to approximate $\pi(a|s)$.

Use policy network to **control** the agent.

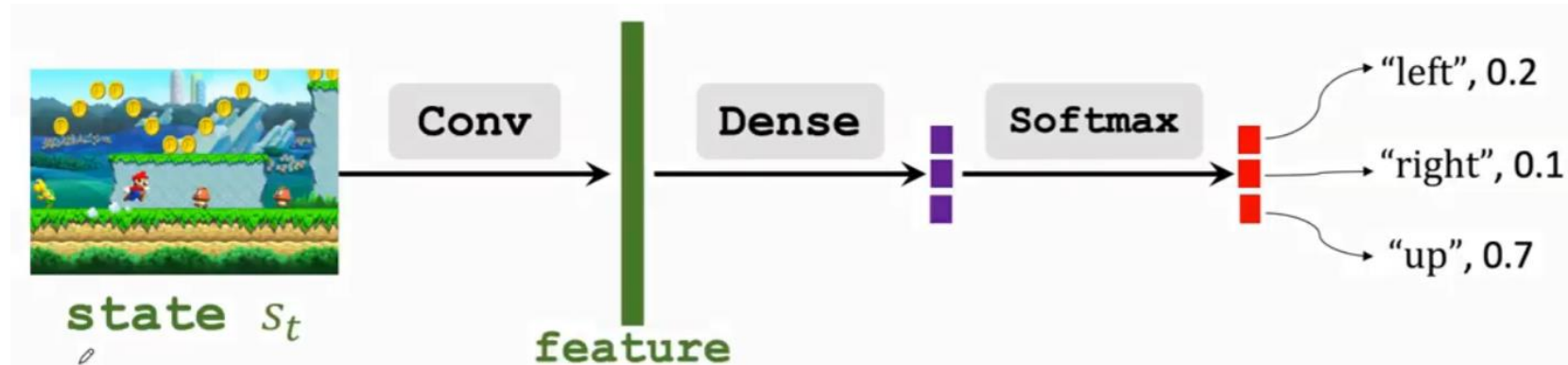
- **Value network – critic**

Use neural network $q(s, a; \mathbf{w})$ to approximate $Q_\pi(s, a)$.

It's only used to give a grade, **not** to control the agent.

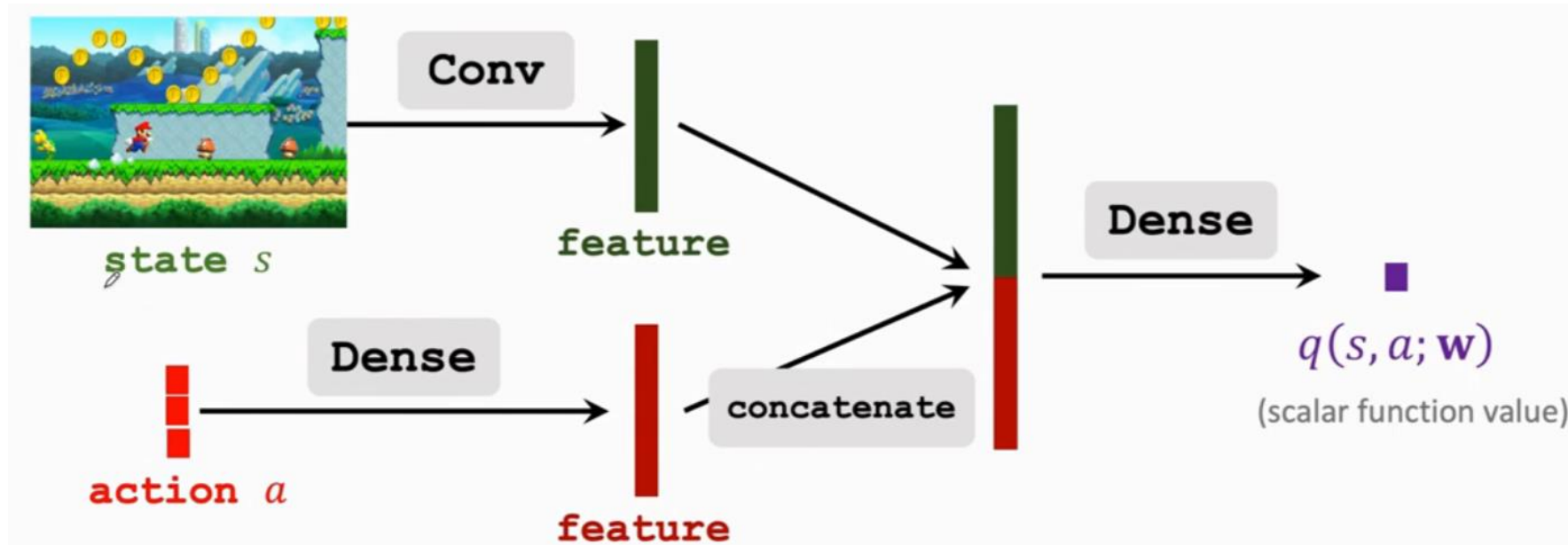
REINFORCE & actor-critic

Policy network – actor: $\pi(a|s; \theta)$



REINFORCE & actor-critic

Value network – critic: $q(s, a; w)$



REINFORCE & actor-critic

Actor-critic method: use neural network to approximate Q_π

- **State-value function**

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

- **Policy network – actor**

$\pi(a|s; \theta)$ is updated to increase state-value $V(s; \theta)$.

- **Value network – critic**

$q(s, a; \mathbf{w})$ is updated to estimate the return more accurately.



REINFORCE & actor-critic

Train actor-critic

1. Observe s_t .
2. Randomly sample 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 \mathbf{w} in value network $q(s, a; \mathbf{w})$, using TD.
5. Update θ in policy network $\pi(a | s; \theta)$, using policy gradient.



REINFORCE & actor-critic

Summary of actor-critic

1. Observe s_t and randomly sample a_t according to $\pi(\cdot | s_t; \theta_t)$.
2. Perform a_t and observe new state s_{t+1} and reward r_t .
3. Randomly sample \tilde{a}_{t+1} according to $\pi(\cdot | s_{t+1}; \theta_t)$. (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: $\delta_t = q_t - (r_t + \gamma \cdot q_{t+1})$. TD target
6. Differentiate value network: $\mathbf{d}_{\mathbf{w},t} = \frac{\partial q(s_t, a_t; \mathbf{w})}{\partial \mathbf{w}} \Big|_{\mathbf{w}=\mathbf{w}_t}$.
7. Update value network $q(s, a; \mathbf{w})$: $\mathbf{w}_{t+1} = \mathbf{w}_t - \alpha \cdot \delta_t \cdot \mathbf{d}_{\mathbf{w},t}$.
8. Differentiate value network: $\mathbf{d}_{\theta,t} = \frac{\partial \log \pi(a_t | s_t; \theta)}{\partial \theta} \Big|_{\theta=\theta_t}$.
9. Update policy network $\pi(a | s; \theta)$: $\theta_{t+1} = \theta_t + \beta \cdot q_t \cdot \mathbf{d}_{\theta,t}$.



REINFORCE & actor-critic

Summary of actor-critic

1. Observe s_t and randomly sample a_t according to $\pi(\cdot | s_t; \theta_t)$.
2. Perform a_t and observe new state s_{t+1} and reward r_t .
3. Randomly sample \tilde{a}_{t+1} according to $\pi(\cdot | s_{t+1}; \theta_t)$. (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: $\delta_t = q_t - (r_t + \gamma \cdot q_{t+1})$. TD target
6. Differentiate value network: $\mathbf{d}_{\mathbf{w},t} = \frac{\partial q(s_t, a_t; \mathbf{w})}{\partial \mathbf{w}} \Big|_{\mathbf{w}=\mathbf{w}_t}$.
7. Update value network $q(s, a; \mathbf{w})$: $\mathbf{w}_{t+1} = \mathbf{w}_t - \alpha \cdot \delta_t \cdot \mathbf{d}_{\mathbf{w},t}$.
8. Differentiate value network: $\mathbf{d}_{\theta,t} = \frac{\partial \log \pi(a_t | s_t; \theta)}{\partial \theta} \Big|_{\theta=\theta_t}$.
9. Update policy network $\pi(a | s; \theta)$: $\theta_{t+1} = \theta_t + \beta \cdot \delta_t \cdot \mathbf{d}_{\theta,t}$.
With baseline



REINFORCE with baseline

State-value function:

$$V_{\pi}(s) = E_A[Q_{\pi}(s, A)] = \sum_a \pi(a|s; \theta) \cdot Q_{\pi}(s, a).$$

Policy gradient:

$$\frac{\partial V_{\pi}(s)}{\partial \theta} = E_{A \sim \pi(\cdot|s; \theta)} \left[\frac{\partial \log \pi(A|s; \theta)}{\partial \theta} \cdot Q_{\pi}(s, A) \right]$$

Baseline: make b independent of A

• Prove $E_{A \sim \pi} \left[b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} \right] = 0$:

$$\begin{aligned} E_{A \sim \pi} \left[b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} \right] &= b \cdot E_{A \sim \pi} \left[\frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} \right] \\ &= b \cdot \sum_a \pi(a|s; \theta) \cdot \left[\frac{\partial \ln \pi(a|s; \theta)}{\partial \theta} \right] \\ &= b \cdot \sum_a \pi(a|s; \theta) \cdot \left[\frac{1}{\pi(a|s; \theta)} \cdot \frac{\partial \pi(a|s; \theta)}{\partial \theta} \right] \\ &= b \cdot \sum_a \frac{\partial \pi(a|s; \theta)}{\partial \theta} \quad \sum_a \pi(a|s; \theta) = 1 \\ &= b \cdot \frac{\partial \left(\sum_a \pi(a|s; \theta) \right)}{\partial \theta} \\ &= b \cdot \frac{\partial 1}{\partial \theta} = 0 \end{aligned}$$

$$\begin{aligned} E_{A \sim \pi(\cdot|s; \theta)} \left[\frac{\partial \log \pi(A|s; \theta)}{\partial \theta} \cdot Q_{\pi}(s, A) \right] &= E_{A \sim \pi(\cdot|s; \theta)} \left[\frac{\partial \log \pi(A|s; \theta)}{\partial \theta} \cdot Q_{\pi}(s, A) \right] - 0 \\ &= E_{A \sim \pi(\cdot|s; \theta)} \left[\frac{\partial \log \pi(A|s; \theta)}{\partial \theta} \cdot Q_{\pi}(s, A) \right] - E_{A \sim \pi} \left[b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} \right] \\ &= E_{A \sim \pi(\cdot|s; \theta)} \left[\frac{\partial \log \pi(A|s; \theta)}{\partial \theta} \cdot (Q_{\pi}(s, A) - b) \right] \end{aligned}$$



REINFORCE with baseline

Policy gradient with baseline

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

Monte carlo approximation

- Randomly sample $a_t \sim \pi(\cdot | s_t; \theta)$ and compute $g(a_t)$
- $g(a_t)$ is an unbiased estimate of the policy gradient:

$$E_{A_t \sim \pi} [g(A_t)] = \frac{\partial V_{\pi}(s_t)}{\partial \theta}.$$

Stochastic policy gradient

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

Stochastic policy gradient ascent:

$$\theta \leftarrow \theta + \beta \cdot g(a_t).$$

b (independent of A_t) **doesn't affect** $E_{A_t \sim \pi} [g(A_t)]$, but **affects** $g(a_t)$. A well chosen b leads to small variance and speeds up convergence.



REINFORCE with baseline

Stochastic policy gradient ($\mathbf{b} = V_\pi(s_t)$)

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

- (Monte carlo) In REINFORCE:

$$u_t \approx Q_\pi(s_t, a_t).$$

- Approximate $V(s; \boldsymbol{\theta})$ with value network $v(s; \mathbf{w})$.

Approximate policy gradient:

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

REINFORCE with baseline

Summary of REINFORCE with baseline

- Play a game to the end and observe 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 $u_t = \sum_{k=t}^n \gamma^{k-t} \cdot r_k$ and $\delta_t = v(s_t; \mathbf{w}) - u_t$.
- Update the policy network:

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

- Update the value network:

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

Advantage actor-critic (A2C)

Stochastic policy gradient ($\mathbf{b} = V_{\pi}(s_t)$)

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

- **Policy network** – actor: $\pi(a | s; \theta)$
 - Approximation to policy function $\pi(a | s)$.
 - Controls the agent.
- **Value network** – critic: $v_{\pi}(s; \mathbf{w})$.
 - Approximation to state value function $V_{\pi}(s)$.
 - Evaluates the state s .



Advantage actor-critic (A2C)

Train A2C

- Observe a transition (s_t, a_t, r_t, s_{t+1}) .
- TD target: $y_t = r_t + \gamma \cdot v(s_{t+1}, \mathbf{w})$.
- TD error: $\delta_t = v(s_t, \mathbf{w}) - y_t$.
- Update policy network (actor) with:

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

- Update the value network (critic) with:

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



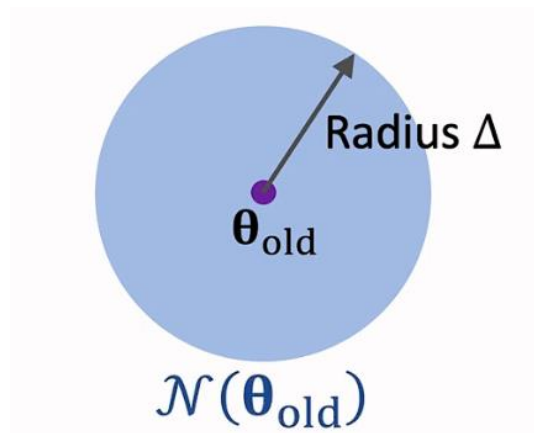
TRPO: trust region policy optimization

Trust region methods, target: $\max_{\theta} J(\theta)$

- Trust region: $N(\theta_{now}) = \{\theta \mid \|\theta - \theta_{now}\|_2 \leq \Delta\}$.
- Construct function $L(\theta \mid \theta_{now})$, satisfying:

$L(\theta \mid \theta_{now})$ is very close to $J(\theta)$, $\forall \theta \in N(\theta_{now})$

- $J(\theta)$ can be replaced by $L(\theta \mid \theta_{now})$ when $\theta \in N(\theta_{now})$



TRPO: trust region policy optimization

Train

- **Approximate**

a) Present policy network parameter is θ_{now} . Use $\pi(a|s; \theta_{now})$ to control the agent, record trajectories: $s_1, a_1, r_1, s_2, a_2, r_2, \dots, s_n, a_n, r_n$.

b) For all t, calculate u_t .

c) Approximate function: $\tilde{L}(\theta|\theta_{now}) = \frac{1}{n} \sum_{t=1}^n \frac{\pi(a_t|s_t;\theta)}{\pi(a_t|s_t;\theta_{now})} \cdot u_t$.

- **Maximize**

$$\theta_{new} = \operatorname{argmax}_{\theta} \tilde{L}(\theta|\theta_{now}); \quad s.t. \|\theta - \theta_{now}\|_2 \leq \Delta.$$

Entropy regularization

- **Entropy**

$$\text{Entropy}(p) = - \sum_{i=1}^n p_i \cdot \ln p_i.$$

- Small entropy means probability is too concentrated, while **big** entropy means greater randomness (**preferred**).

- **Entropy of probability distribution**

$$H(s; \theta) \triangleq \text{Entropy}[\pi(\cdot | s; \theta)] = - \sum_{a \in A} \pi(a | s; \theta) \cdot \ln \pi(a | s; \theta).$$

- **Policy-based RL using entropy regularization**

$$\max_{\theta} J(\theta) + \lambda \cdot E_S[H(s; \theta)].$$

- **Gradient and update θ**

$$\tilde{g}(s, a; \theta) \triangleq [Q_{\pi}(s, a) - \lambda \cdot \ln \pi(a | s; \theta) - \lambda] \cdot \nabla_{\theta} \ln \pi(a | s; \theta).$$

$$\theta \leftarrow \theta + \beta \cdot \tilde{g}(s, a; \theta).$$

RNN: recurrent neural network

- **Problem:** in many cases, the observations are partial, leading to difficult decision-making. Replacing s_t with o_t is the simplest way, but it's not efficient enough.
- **Solution: remember all the o_i :** o_1, o_2, \dots, o_t , denoted by $o_{1:t} = [o_1, o_2, \dots, o_t]$, and policy network is $\pi(a_t | o_{1:t}; \theta)$.
- However, with the growth of t , the length of $o_{1:t}$ changes accordingly.
- We can use recurrent neural network (RNN).

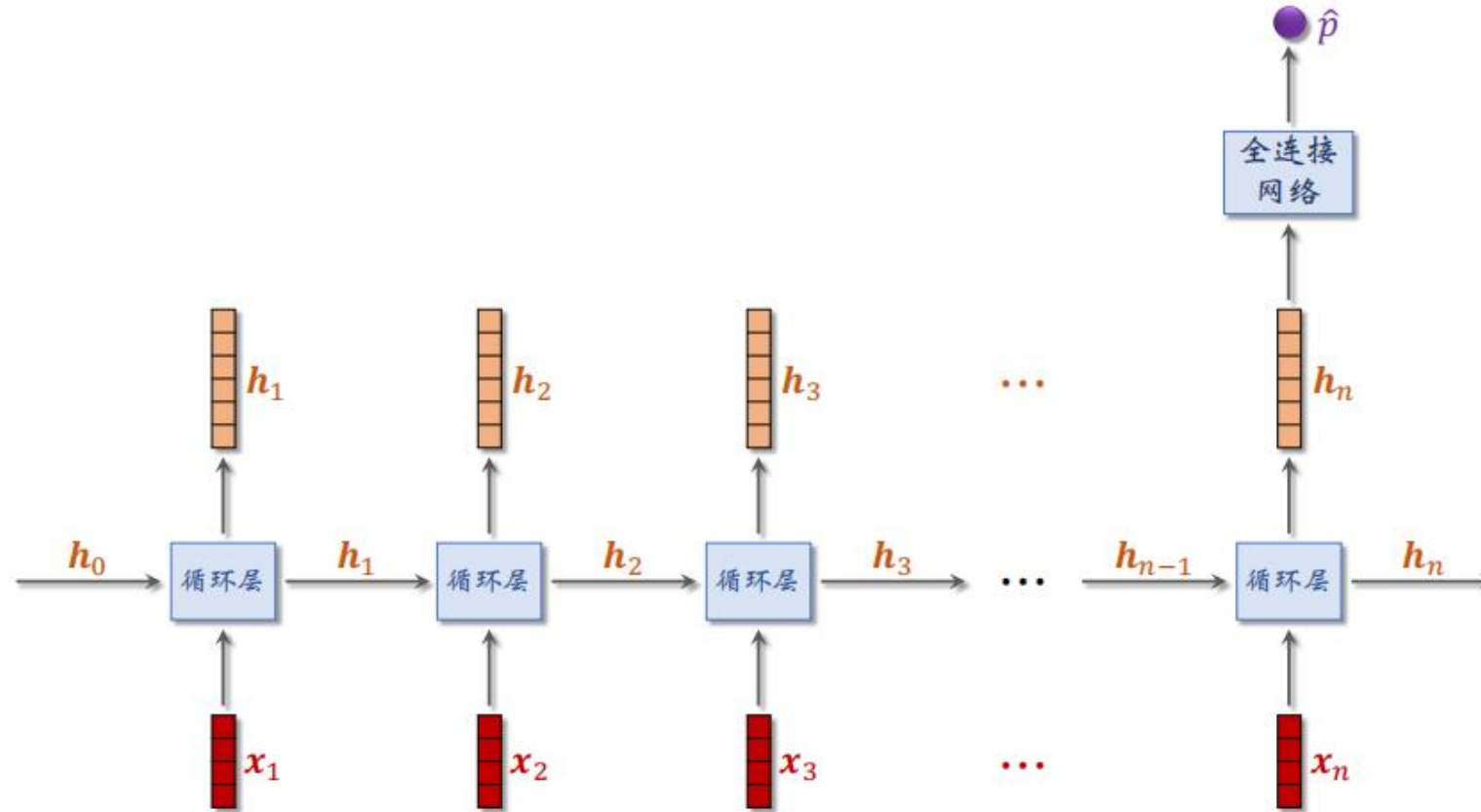
RNN: recurrent neural network

- For all $t=1, 2, \dots, n$, recurrent layer maps $[x_1, x_2, \dots, x_t]$ to h_t .

$$\begin{array}{ll}
 (\mathbf{x}_1) & \implies \mathbf{h}_1, \\
 (\mathbf{x}_1, \mathbf{x}_2) & \implies \mathbf{h}_2, \\
 (\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3) & \implies \mathbf{h}_3, \\
 & \vdots \\
 (\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \dots, \mathbf{x}_{n-1}) & \implies \mathbf{h}_{n-1}, \\
 (\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \dots, \mathbf{x}_{n-1}, \mathbf{x}_n) & \implies \mathbf{h}_n.
 \end{array}$$

- The length of h_t is fixed.

RNN: recurrent neural network

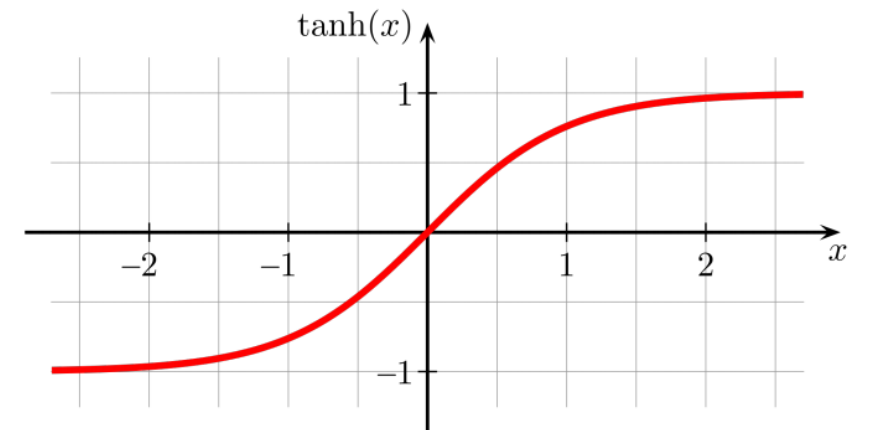
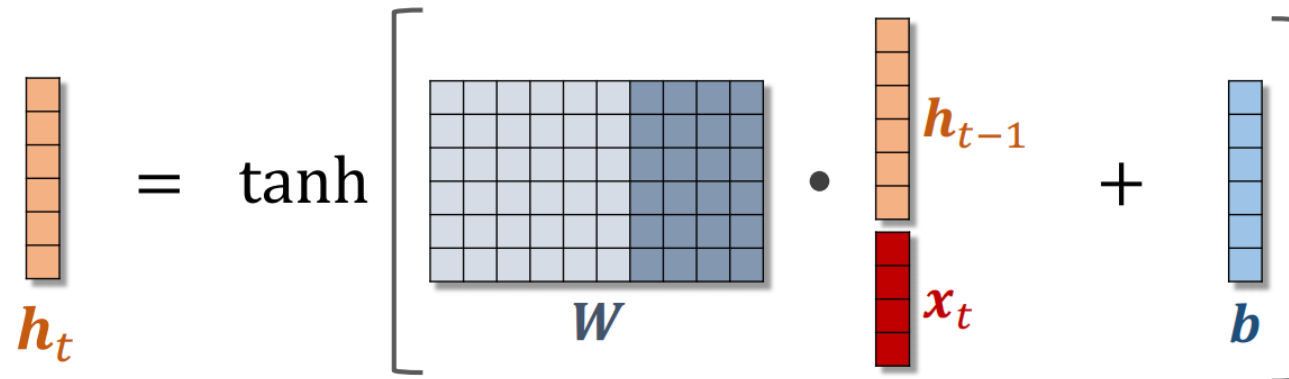


Input is $[x_1, x_2, \dots, x_n]$, output is \hat{p} .

RNN: recurrent neural network

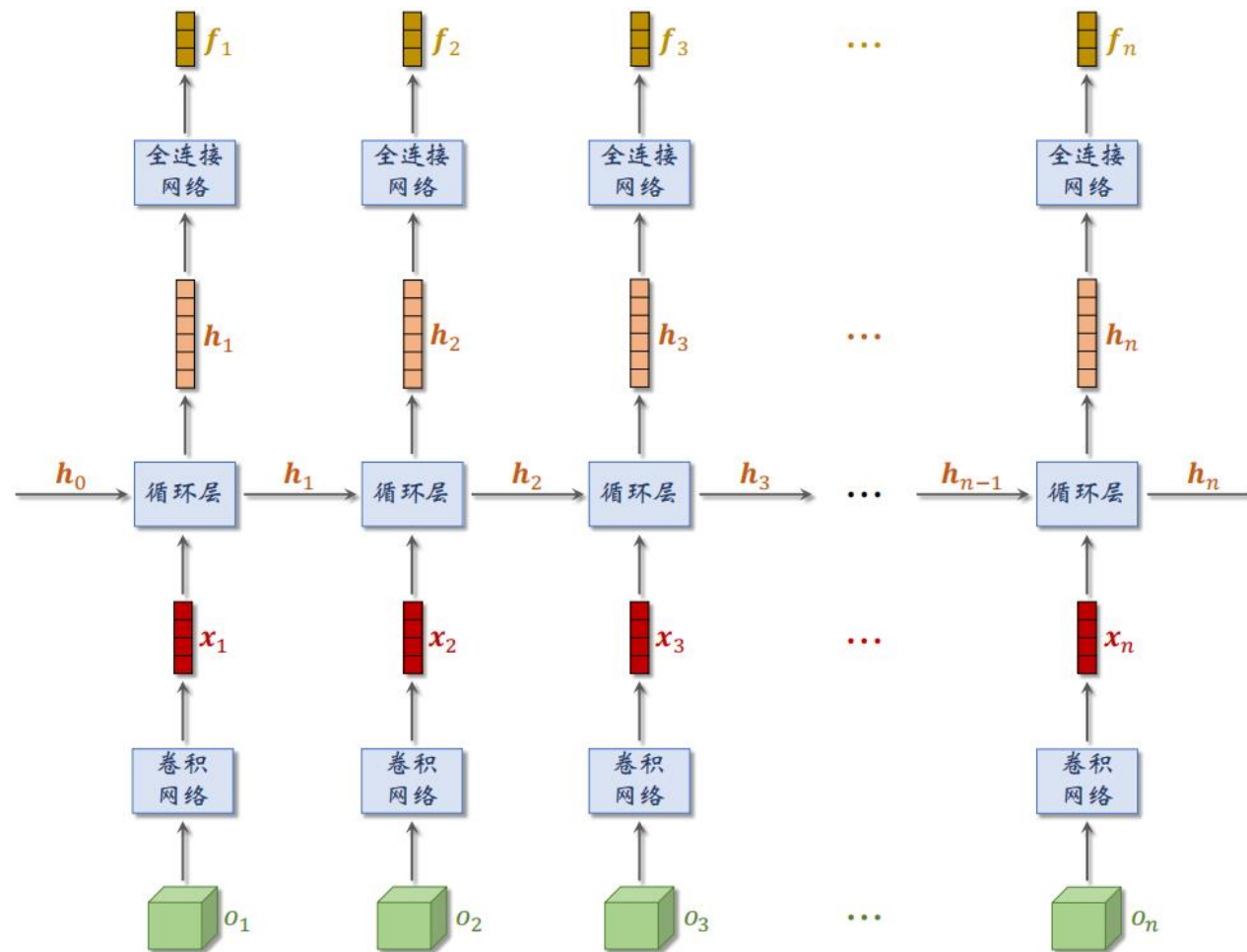
- Simple recurrent layer

$$h_t = \tanh(W[h_{t-1}; x_t] + b).$$



RNN: recurrent neural network

RNN as policy network: $f_t = \pi(a_t | o_{1:t}; \theta)$.



Thank you.

