

# RL using Function Approximation\_ For continuous space

---

## RL using Function Approximation\_ For continuous space

### 1. Q-function Approximation

Background & Basic Idea

Method 1: use linearly parametrized approximation

Q-function Approximation

Action Space Discretization

Fuzzy Q-iteration

Model

Policy

Fuzzy Q-iteration

### 2. Actor-critic methods

Structure

Update Critic: Value Estimation

Update Actor: Policy Update

### 3. One Really Classical Example

## 1. Q-function Approximation

---

### Background & Basic Idea

In real-life control,  $X$ ,  $U$  continuous

- approximate Q-function  $\hat{Q}$  must be used

### Method 1: use linearly parametrized approximation

## Q-function Approximation

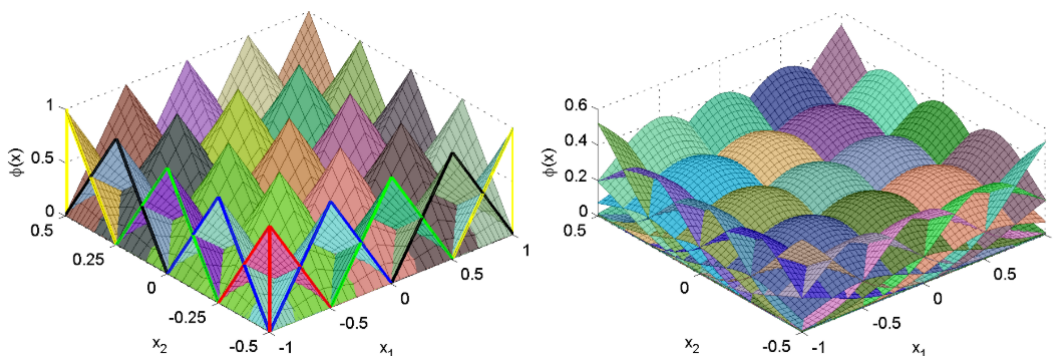
- use **Basis Function** to approximate Q-function in a continuous space

$$\hat{Q} = \sum_{i=1}^N \theta_i \phi_i(x, u)$$

$$\phi_i(x, u) : X \times U \mapsto \mathbb{R}$$

usually normalized:  $\sum_i \phi_i(x) = 1$

- E.g., fuzzy approximation, RBF network approximation



- Policy is greedy in  $\hat{Q}$ , computed on demand for given  $x$ :

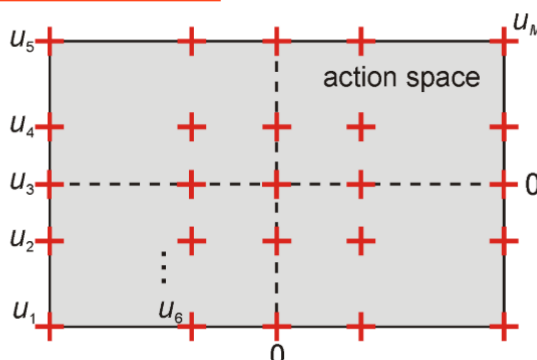
$$\pi(x) = \arg \max_u \hat{Q}(x, u)$$

- Approximator must ensure efficient arg max solution

## Action Space Discretization

- Choose  $M$  discrete actions  $u_1, \dots, u_M \in U$
- Solve “arg max” by explicit enumeration
- In a control problem, we always make the grid more detail around the attractor, for example, 0 in the graph

Example: grid discretization



# Fuzzy Q-iteration

## Model

Given:

- N basis functions  $\phi_1, \dots, \phi_N$
- M discrete actions  $u_1, \dots, u_M$

Store

- N × M matrix of parameters  $\theta$  (one for each pair basis function–discrete action)
  - same row: same  $\phi$
  - same column: same action

$$\hat{Q}^\theta(x, u_j) = \sum_{i=1}^N \phi_i(x) \theta_{i,j} = [\phi_1(x) \dots \phi_N(x)] \begin{bmatrix} \theta_{1,j} \\ \vdots \\ \theta_{N,j} \end{bmatrix}$$

## Policy

$$\hat{\pi}^*(x) = \arg \max_{u_j, j=1, \dots, M} \hat{Q}^{\theta^*}(x, u_j)$$

( $\theta^*$  = converged parameter matrix )

## Fuzzy Q-iteration

### Fuzzy Q-iteration

**repeat** at each iteration  $\ell$

**for all** cores  $x_i$ , discrete actions  $u_j$  **do**

$$\theta_{\ell+1,i,j} = \rho(x_i, u_j) + \gamma \max_{j'} \hat{Q}^{\theta_\ell}(f(x_i, u_j), u_{j'})$$

**end for**

**until** convergence

↘ weighted sum

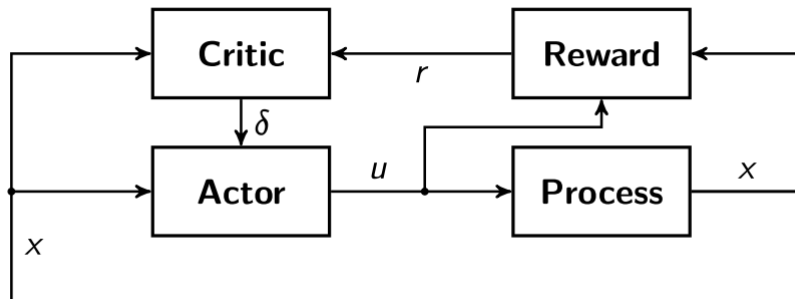
## 2. Actor-critic methods

---

# Structure

## Explicitly separated value function and policy

- Actor = control policy  $\pi(x)$
- Critic = state value function  $V(x)$



### Actor-critic

**for** every trial **do**

    initialize  $x_0$ , choose initial action  $u_0 = \tilde{u}_0$

**repeat** for each step  $k$

        apply  $u_k$ , measure  $x_{k+1}$ , receive  $r_{k+1}$

        choose **next** action  $u_{k+1} = \hat{\pi}(x_{k+1}, \varphi_k) + \tilde{u}_{k+1}$

$\Delta_k = r_{k+1} + \hat{V}(x_{k+1}, \theta_k) - \hat{V}(x_k, \theta_k)$

$\theta_{k+1} = \theta_k + \alpha_c \Delta_k \left. \frac{\partial \hat{V}(x, \theta)}{\partial \theta} \right|_{x=x_k, \theta=\theta_k}$

$\varphi_{k+1} = \varphi_k + \alpha_a \Delta_k \tilde{u}_k \left. \frac{\partial \hat{\pi}(x, \varphi)}{\partial \varphi} \right|_{x=x_k, \varphi=\varphi_k}$

**until** terminal state

**end for**

### When facing continuity:

- Actor parameterized in  $\phi : \hat{\pi}(x; \phi)$
- Critic parameterized in  $\theta : \hat{V}(x; \theta)$

Parameters  $\phi$  and  $\theta$ , have finite size, but approximate functions on continuous (infinitely large) spaces

## Update Critic: Value Estimation

The task of the critic is to **predict the expected future reinforcement  $r$  the process will receive being in the current state and following the current control policy.**

- For doing that, we need to train critic, the prediction error is always a train input:

Use sample  $(x_k; u_k; x_{k+1}; r_{k+1})$  at each step k and parameterized V:

Note that both  $\hat{V}(s_k)$  and  $\hat{V}(s_{k+1})$  are known at time k, since  $\hat{V}(s_{k+1})$  is a prediction obtained for the current process state: For example, assuming a NN, we can always get a prediction of next state value based on current state, we get the  $\Delta$  from **real system feedback**  $r_{k+1}$  and **critic original prediction**

$$\Delta_k = V(s_k) - \hat{V}(s_k) = r_{k+1} + \gamma \hat{V}(s_{k+1}) - \hat{V}(s_k)$$

- For example, Let the critic be represented by a neural network or a fuzzy system, in my opinion, **the index of  $\theta$  represents the iteration version, it accidentally equal to state/step number, because each step we make a update. It does not mean the corresponding  $\theta$  of state k or k+1**

$$\hat{V}(s_{k+1}) = \hat{V}(s_{k+1}, \theta_k)$$

To update  $\theta_k$ , a gradient-descent learning rule is used:

$$\theta_{k+1} = \theta_k + \alpha_c \Delta_k \frac{\partial \hat{V}(s_k, \theta_k)}{\partial \theta_k}$$

$\alpha_c > 0$ , learning rate of critic

$$\Delta_k > 0, \text{ i.e., } r_{k+1} + \gamma \hat{V}^\pi(x_{k+1}, \theta_k) > \hat{V}^\pi(x_k, \theta_k)$$

$\Rightarrow$  old estimate too low, increase  $\hat{V}$

$$\Delta_k < 0, \text{ i.e., } r_{k+1} + \gamma \hat{V}^\pi(x_{k+1}, \theta_k) < \hat{V}^\pi(x_k, \theta_k)$$

$\Rightarrow$  old estimate too high, decrease  $\hat{V}$

## Update Actor: Policy Update

The actor (i.e., the policy) can be adapted in order to **establish an optimal mapping between the system states and the control actions.**

$$u_k = \hat{\pi}(x_k, \varphi_k) + \tilde{u}_k, \hat{\pi} = \text{actor}, \tilde{u}_k = \text{exploration}$$

$$\varphi_{k+1} = \varphi_k + \alpha_a \Delta_k \tilde{u}_k \left. \frac{\partial \hat{\pi}(x, \varphi)}{\partial \varphi} \right|_{x=x_k}$$

$$\Delta_k > 0, \text{ i.e., } r_{k+1} + \gamma \hat{V}^\pi(x_{k+1}, \theta_k) > \hat{V}^\pi(x_k, \theta_k)$$

$\Rightarrow \tilde{u}_k$  had positive effect, move in that direction

$\Delta_k < 0$ , i.e.,  $r_{k+1} + \gamma \hat{V}^\pi(x_{k+1}, \theta_k) < \hat{V}^\pi(x_k, \theta_k)$   
 $\Rightarrow \tilde{u}_k$  had negative effect, move away from that direction

### 3. One Really Classical Example

---

Lecture Notes, Application 1