



# システム制御理論と統計的機械学習

## 第 9 章：強化学習

---

加嶋 健司

October 10, 2025

京都大学 情報学研究科

# 本章の流れ

## 9.1 最適制御と強化学習

強化学習の概要

## 9.2 モデルベース設計

方策反復法と方策勾配法

## 9.3 モデルフリー設計

行動価値関数をデータから推定

# 最適制御と強化学習

---

# 最適制御問題

基本的には以下の無限時間最適制御問題を解くことが目標

## Problem (時不変無限区間最適制御問題)

遷移確率密度関数  $\Psi$  と初期状態  $x_0 \in \text{rv}(\mathbb{X})$  をもつ確率システム,

非負のステージコスト  $\ell : \mathbb{X} \times \mathbb{U} \rightarrow \mathbb{R}_+$ , 集合値関数  $U(x) \subset \mathbb{U}$ ,  $x \in \mathbb{X}$  および  $\beta \in (0, 1)$

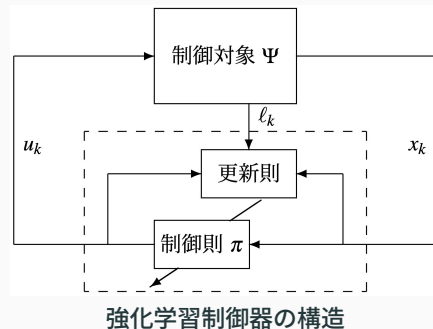
$\pi(x_k) \in U(x_k)$  を満たし,

$$J(\pi) := \mathbb{E}^\pi \left[ \sum_{k=0}^{\infty} \beta^k \ell(x_k, u_k) \right] \quad (9.1)$$

を最小化する制御則  $\pi$  を求めよ.

# モデルフリー設計

- モデルベース設計：遷移確率密度関数  $\Psi$  を利用する設計法（たとえば  $\Psi$  で定まるベルマン方程式を解く）
- 強化学習（reinforcement learning）： $\Psi$  は未知であるとし（モデルフリー設計）， $(x_k, u_k, \ell_k)$  のデータを用いて制御則  $\pi$  を更新して制御器を設計すること
- 制御則を  $\pi^*$  に収束させる更新則を設計  
更新に有用な情報を得るための入力設計



# モデルベースト設計

---

# ベルマン方程式 (1) 価値関数

## 作用素

$$\mathcal{A}_x[V](x, u) := \int_{\mathbb{X}} V(x') \Psi(x'|x, u) dx' = \mathbb{E}[V(x_{k+1}) | x_k = x, u_k = u] \quad (9.2)$$

$$\mathcal{B}[V](x) := \inf_{u \in U(x)} (\ell + \beta \mathcal{A}_x[V])(x, u) \quad (9.3)$$

$$\mathcal{B}^\pi[V](x) := \int_{U(x)} (\ell + \beta \mathcal{A}_x[V])(x, u) \pi(u|x) du \quad (9.4)$$

## 状態価値関数の定義

$$V^\pi(x) := \mathbb{E}_{x_0=x}^\pi \left[ \sum_{k=0}^{\infty} \beta^k \ell_k \right] \quad (9.5)$$

$$V^*(x) := \inf_{\pi} V^\pi(x) \quad (9.6)$$

## 行動価値関数の定義

$$Q^\pi(x, u) := \mathbb{E}_{x_0=x, u_0=u}^\pi \left[ \sum_{k=0}^{\infty} \beta^k \ell_k \right] \quad (9.7)$$

$$Q^*(x, u) := \inf_{\pi} Q^\pi(x, u) \quad (9.8) \quad 4$$

## ベルマン方程式 (2) 再掲

$$Q^*(x, u) = (\ell + \beta \mathcal{A}_x[V^*])(x, u)$$

$$V^*(x) = \mathcal{B}[V^*](x) = \inf_{u \in U(x)} Q^*(x, u)$$

$$Q^\pi(x, u) = (\ell + \beta \mathcal{A}_x[V^\pi])(x, u)$$

$$V^\pi(x) = \mathcal{B}^\pi[V^\pi](x) = \int_{U(x)} Q^\pi(x, u) \pi(u|x) du$$

$$\pi^*(x) = \arg \min_{u \in U(x)} Q^*(x, u)$$

状態価値関数・行動価値関数の入力

	$x_0$	$u_0$	$u_k, k \geq 1$
$V^*(x)$	$x$	$\pi^*(x)$	$\pi^*(x_k)$
$Q^*(x, u)$	$x$	$u$	$\pi^*(x_k)$
$V^\pi(x)$	$x$	$\pi(x)$	$\pi(x_k)$
$Q^\pi(x, u)$	$x$	$u$	$\pi(x_k)$



## ベルマン方程式 (3) 価値反復法

### 補題 9.2.1 – $\mathcal{B}^\pi$ の単調性

関数  $V_1, V_2 : \mathbb{X} \rightarrow \mathbb{R}$  が  $V_1 \leq V_2$  を満たすならば,  $\mathcal{B}^\pi[V_1] \leq \mathcal{B}^\pi[V_2]$  である.

### 補題 9.2.2 – 終端コストと状態価値関数

$$(\mathcal{B}^\pi)^{\bar{k}}[V](x) = \mathbb{E}_{x_0=x}^\pi \left[ \sum_{k=0}^{\bar{k}-1} \beta^k \ell(x_k, u_k) + \beta^{\bar{k}} V(x_{\bar{k}}) \right] \quad (9.14)$$

### 定理 9.2.3 – $V^\pi$ に対する価値反復法 ( $\ell, V$ は有界)

$$\lim_{k \rightarrow \infty} \|(\mathcal{B}^\pi)^k[V] - V^\pi\|_\infty = 0 \quad (9.15)$$

## ベルマン方程式 (4) 方策反復法

**方策反復法 (policy iteration) :**  $\pi_+(x) := \arg \min_{u \in U(x)} Q^\pi(x, u)$  により制御則を更新

### 定理 9.2.4 – 方策改善定理 (policy improvement theorem)

$\ell$  は有界な関数,  $\pi$  を任意の状態フィードバック制御則とする. このとき,

$$Q^\pi(x, \pi_+(x)) \leq V^\pi(x), \forall x \in \mathbb{X} \quad (9.17)$$

を満たす任意の  $\pi_+$  に対して,  $V^{\pi_+} \leq V^\pi$  が成り立つ. 逆に,  $V^\pi$  が有界かつ

$$Q^\pi(x, u) \geq V^\pi(x) - \varepsilon(1 - \beta), \forall x \in \mathbb{X}, u \in U(x) \quad (9.18)$$

が成り立つとき,  $V^\pi \leq V^* + \varepsilon$  である.

- 初期時刻以外は  $\pi$  で固定したままで, 初期時刻の入力のみ変更して性能が改善しなければ,  $\pi$  は既に最適である (初期時刻以外を同時に探索しても改善しない)

## ベルマン方程式 (5) LQR に対する方策反復法

LQR に対する方策反復法は、リカッチ方程式の反復解法 (Hewer, 1971)

$$x_{k+1} = Ax_k + Bu_k + v_k, \ell_k = \begin{bmatrix} x_k^\top & u_k^\top \end{bmatrix} \begin{bmatrix} Q & S \\ S^\top & R \end{bmatrix} \begin{bmatrix} x_k \\ u_k \end{bmatrix} \quad (9.21)$$

1. 制御則  $\pi$  を  $u_k = -Kx_k$  としたとき, 状態価値関数  $V^\pi(x) = x^\top \Pi x + \gamma$  ( $\because$  定理 6.2.2)

$$\beta(A - BK)^\top \Pi (A - BK) + \begin{bmatrix} I & -K^\top \end{bmatrix} \begin{bmatrix} Q & S \\ S^\top & R \end{bmatrix} \begin{bmatrix} I \\ -K \end{bmatrix} = \Pi \quad (9.22)$$

2. 行動価値関数は, ( $\because Q^\pi(x, u) = (\ell + \beta \mathcal{A}_x[V^\pi])(x, u)$ )

$$Q^\pi(x, u) = \begin{bmatrix} x^\top & u^\top \end{bmatrix} \Upsilon^\pi \begin{bmatrix} x \\ u \end{bmatrix} + \gamma_\pi, \Upsilon^\pi := \begin{bmatrix} * & S + \beta A^\top \Pi B \\ * & R + \beta B^\top \Pi B \end{bmatrix} \quad (9.23)$$

3. 更新則は,

$$\arg \min_{u \in \mathbb{R}^{n_u}} Q^\pi(x, u) = -K_+ x := -(\Upsilon^\pi)_{22}^{-1} (\Upsilon^\pi)_{12}^\top x \quad (9.24)$$

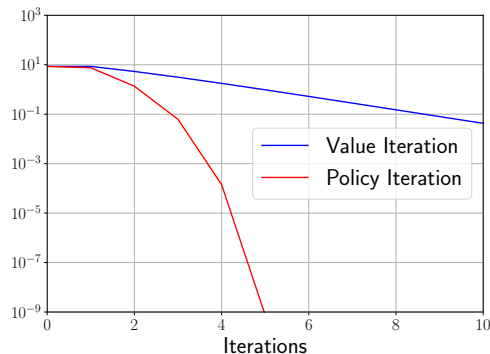
## 方策勾配法 (6) LQR に対する反復法の例

- ・ 価値反復法であるリカッチイタレーションでは、特別な初期化は不要であるが、反復中に得られる制御則は制御区間を有限で打ち切った問題の最適制御性能しか考慮しておらず、十分収束するまでは安定性も保証されない
- ・ 方策反復法は安定化制御則で初期化する必要があるが、反復ごとに無限時間の制御性能が改善される

### 例 9.2.5 – 方策反復法

方策反復法では各反復においてリアプノフ方程式を解く必要があるが、誤差の収束は価値反復法よりも速い。

- ・ 価値反復法  $\Pi_{i+1} = \text{Ric}(\Pi_i)$



## 方策勾配法 (1) 勾配の計算

### 方策勾配法 (policy gradient method) :

制御則をパラメータ  $\theta \in \mathbb{R}^{n_p}$  で規定し、コストの  $\theta$  に関する勾配を用いて最適化

- ある  $\theta$  に対する  $Q^{\pi_\theta}$  が既知でも、勾配計算は自明でない
- 確定方策  $\pi_\theta(x)$  のとき,

$$J(\pi_\theta) = \mathbb{E}[V^{\pi_\theta}(x_0)] = \mathbb{E}[Q^{\pi_\theta}(x_0, \pi_\theta(x_0))] \quad (9.25)$$

の  $\theta$  に関する勾配をとればよいが、 $\pi_\theta$  はこちらが自由に設定できる関数であるものの、 $Q^{\pi_\theta}$  の関数形の  $\theta$  への依存性は自明ではない。

## 方策勾配法 (2) 方策勾配定理

### 定理 9.2.6 – 方策勾配定理

確率方策  $\pi_\theta(u|x)$

$$\nabla_\theta J(\pi_\theta) = \mathbb{E}^{\pi_\theta} \left[ \sum_{k=0}^{\infty} \beta^k \int_{U(x_k)} (\nabla_\theta \pi_\theta(u|x_k)) Q^{\pi_\theta}(x_k, u) du \right] \quad (9.26)$$

確定方策  $\pi_\theta(x)$

$$\nabla_\theta J(\pi_\theta) = \mathbb{E}^{\pi_\theta} \left[ \sum_{k=0}^{\infty} \beta^k \nabla_\theta \pi_\theta(x_k) \nabla_u Q^{\pi_\theta}(x_k, \pi_\theta(x_k)) \right] \quad (9.27)$$

- $Q^{\pi_\theta}$  の  $\theta$  に関する勾配  $\nabla_\theta Q^{\pi_\theta}$  が現れていない
- 期待値が  $\varphi_{x_0}$  ではなく  $\sum_{k=0}^{\infty} \beta^k \varphi_{x_k}^{\pi_\theta}$  により重み付けされた積分

## 方策勾配法 (3) 行動価値関数の必要情報

### 注意 9.2.7 – 各手法に必要な $Q^{\pi_\theta}$ の情報

- 方策反復法の更新則

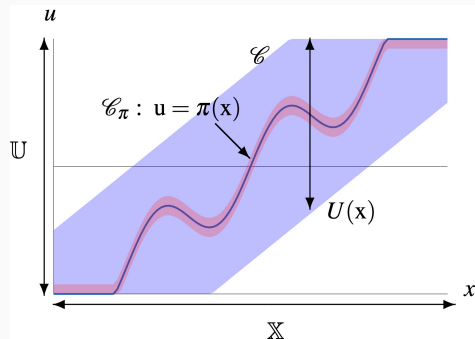
$$\pi_+(x) := \operatorname{argmin}_{u \in U(x)} Q^\pi(x, u)$$

の計算には,  $\mathcal{C} := \{(x, u) : x \in \mathbb{X}, u \in U(x)\}$  全体

- 方策勾配法の勾配計算

$$\mathbb{E}^{\pi_\theta} \left[ \sum_{k=0}^{\infty} \beta^k \nabla_\theta \pi_\theta(x_k) \nabla_u Q^{\pi_\theta}(x_k, \pi_\theta(x_k)) \right]$$

には,  $\mathcal{C}_{\pi_\theta} := \{(x, \pi_\theta(x)) : x \in \mathbb{X}\}$  近傍



行動価値関数を求める領域

# モデルフリー設計

---



# 行動価値関数の推定 (1)

遷移確率密度関数  $\Psi$  を用いずに方策反復法や方策勾配法を実装

## 問題 9.3.1 – 行動価値関数の推定

遷移確率密度関数  $\Psi$  および状態フィードバック制御則  $\pi$  に対して,  $x_:$ ,  $u_:$ ,  $\ell_:$  の実現値系列を用いて, 行動価値関数  $Q^\pi$  を推定せよ.

### アルゴリズム 4 方策反復法

input  $\pi$

output  $\pi_+$

- 1:  $Q^\pi$  を求める
- 2: **for all**  $x$  in  $\mathbb{X}$  **do**
- 3:  $\pi_+(x) \leftarrow \arg \min_{u \in U(x)} Q^\pi(x, u)$

### アルゴリズム 6 方策勾配法

input  $\theta, a$

output  $\theta_+$

- 1:  $Q^{\pi_\theta}$  を求める
- 2:  $\nabla_\theta J(\pi_\theta)$  を求める
- 3:  $\theta_+ \leftarrow \theta - a \nabla_\theta J(\pi_\theta)$

## 行動価値関数の推定 (2) TD 誤差

### TD 誤差 (temporal difference error)

$$\delta_k := \ell_k + \beta \hat{Q}^\pi(x_{k+1}, \pi(x_{k+1})) - \hat{Q}^\pi(x_k, u_k) \quad (9.28)$$

- $Q^\pi(x_k, u_k)$  を  $\ell(x_k, u_k) + \beta Q^\pi(x_{k+1}, \pi(x_{k+1}))$  の近似とする際の誤差
- ベルマン方程式  $Q^\pi(x, u) = (\ell + \beta \mathcal{A}_x[V^\pi])(x, u)$  より,  $\hat{Q}^\pi = Q^\pi$  ならば

$$\mathbb{E}[\delta_k | x_{:k}, u_{:k}, l_{:k-1}] = 0, \quad \forall k, x_{:k}, u_{:k}, l_{:k-1} \quad (9.29)$$

$\delta_k$  が 0 に近い値をとる  $\hat{Q}^\pi$  を  $Q^\pi$  の推定結果とするという立場をとる

## 行動価値関数の推定 (3) TD 学習

$$\delta_k := \ell_k + \beta \hat{Q}^\pi(x_{k+1}, \pi(x_{k+1})) - \hat{Q}^\pi(x_k, u_k)$$

各組  $(x_k, u_k, \ell_k, x_{k+1})$  に対する一段先予測誤差の和の最小化

$$\sum_k \left( \ell_k + \beta \hat{Q}^\pi(x_{k+1}, \pi(x_{k+1})) - \hat{Q}^\pi(x_k, u_k) \right)^2 \quad (9.30)$$

- $\pi(x_{k+1})$  を  $u_{k+1}$  としたり,  $u_k$  を  $\pi(x_k)$  としたりしないように注意
- $u_k = \pi(x_k)$  のもとでの閉ループデータである必要はない (過去に  $\pi$  とは無関係に取得・蓄積されたデータも利用可能)

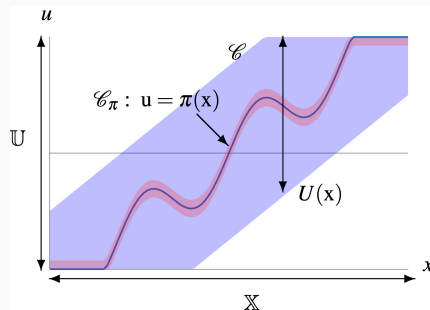
## 行動価値関数の推定 (6) TD 学習の例

確定的な状態フィードバック  $u_k = \pi(x_k)$  のもとでのデータのみでは、 $\mathcal{C}_\pi$  付近の行動価値関数しか推定できない

- 雑音  $z_k$  を付加して

$$u_k = \pi(x_k) + z_k \quad (9.32)$$

- $\varepsilon > 0$  に対して、 $u_k$  を確率  $(1 - \varepsilon)$  で  $\pi(x_k)$ 、確率  $\varepsilon$  で一様乱数とする ( $\varepsilon$ -greedy)
- $T > 0$  に対して、 $\pi(\cdot|x) \propto \exp(-Q(x, \cdot)/T)$  とする (ボルツマン選択則)



行動価値関数を求める領域

## 行動価値関数の推定 (5) RLS フィルタの利用

特徴写像  $\phi : \mathbb{X} \times \mathbb{U} \rightarrow \mathbb{R}_f^n$  を用いて,

$$\hat{Q}^\pi(x, u; \hat{p}) := \hat{p}^\top \phi(x, u), \quad \hat{p} \in \mathbb{R}^{n_f} \quad (9.33)$$

- RLS フィルタを用いて最小化

$$\sum_k \left( l_k - \hat{p}^\top \{ \phi(x_k, u_k) - \beta \phi(x_{k+1}, \pi(x_{k+1})) \} \right)^2 \quad (9.34)$$

- 制御対象が線形るとき,  $Q^\pi$  は  $\begin{bmatrix} x^\top & u^\top \end{bmatrix}^\top$  の二次形式  
 $x, u$  の要素の二次単項式および定数 1 を並べたベクトル  $\phi(x, u)$  を用いて,

$$Q^\pi(x, u) = \begin{bmatrix} x^\top & u^\top \end{bmatrix} \Upsilon^\pi \begin{bmatrix} x \\ u \end{bmatrix} + \gamma_\pi = p^\top \phi(x, u) \quad (9.35)$$

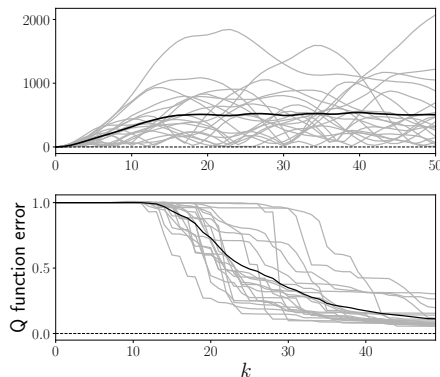
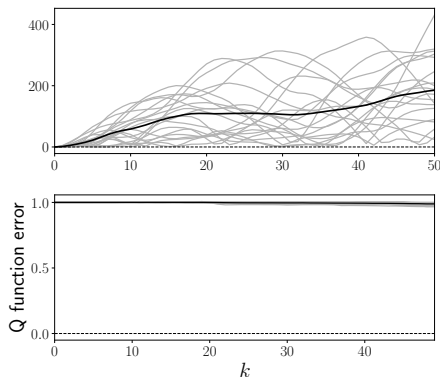
$p$  は  $\Upsilon^\pi$  の対称性から  $(n_x + n_u)(n_x + n_u + 1)/2 + 1$  次元

## 行動価値関数の推定 (6) TD 学習の例

### 例 9.3.3 – TD 学習における探索と利用のトレードオフ

線形システムに対して  $u_k = -K_0 x_k + z_k$  ( $\rho(A - BK_0) = 0.88$ ) で制御された際の軌道

上段：状態の変動，下段：学習誤差，左： $z_k \sim \mathcal{N}(0, 2^2)$ ，右： $z_k \sim \mathcal{N}(0, 10^2)$



# アルゴリズム例 (1) TD 学習にもとづく方策反復法

## ボルツマン選択則を用いる方策反復法と TD 学習の組み合わせ

### アルゴリズム 7 TD 学習を用いた方策反復法

**input**  $k, \pi, l, x, \mathcal{D}$  ▷  $k, \pi_{k-1}, l_{k-1}, x_k$   
**output**  $\pi, u, \mathcal{D}$  ▷  $\pi_k, u_k$   
1: **if**  $k \% k_{\text{update}} = 0$  **then** ▷ 時刻  $k_{\text{update}}$  ごとに制御則を更新  
2:      $\mathcal{D} \cup \{(k, l, x)\}$  に対して式 (9.30) を最小化する  $\hat{Q}^\pi$  を求める ▷  $Q^\pi$  の推定  
3:     **for all**  $x'$  in  $\mathbb{X}$  **do**  
4:          $\pi(\cdot|x') \leftarrow c \exp(-\hat{Q}^\pi(x', \cdot)/T)$  ▷ 式 (9.16) とボルツマン選択則  
5:     **Initialize:**  $\mathcal{D}$  ▷ 必要に応じてデータを初期化  
6: **else**  
7:      $\pi \leftarrow \hat{\pi}$   
8:      $u \sim \pi(\cdot|x)$  ▷ 制御入力を計算  
9:      $\mathcal{D} \leftarrow \mathcal{D} \cup \{(k, l, x, u)\}$  ▷ データを保存

$k_{\text{update}}$  ステップの学習ごとに制御則を更新

## アルゴリズム例 (2) TD 学習にもとづく方策反復法の例

### 例 9.3.4 – TD 学習と方策反復法にもとづく強化学習

- $z_k \sim \mathcal{N}(0, 10^2)$ ,  $k_{\text{update}} = 50$
- 制御則の更新ごとに初期化

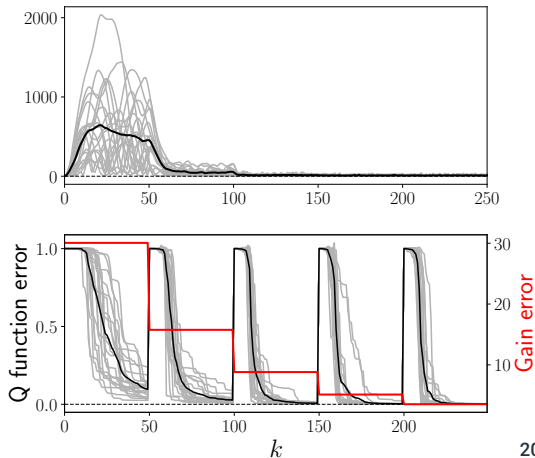
$$\hat{p}_k = 0, \Sigma_k = 10^4 \times I$$

上段：状態の変動

下段：相対近似誤差

Q 関数  $\|\hat{\Upsilon}_k - \Upsilon^{\pi_i}\| / \|\Upsilon^{\pi_i}\|$

ゲイン  $\|K_k - K^*\| / \|K^*\|$





## アルゴリズム例 (3) 制御則の更新間隔

$k_{\text{update}}$  が長いと、行動価値関数の推定精度は向上するが、その間の制御性能は妥協

$$\sum_k \left( l_k + \beta \hat{Q}^\pi(x_{k+1}, \pi(x_{k+1})) - \hat{Q}^\pi(x_k, u_k) \right)^2 \quad (9.30)$$

- $\mathbb{X}, \mathbb{U}$  が有限集合であるマルコフ決定過程の場合、特定の点における  $\hat{Q}^\pi$  の関数値を容易に更新できる ( $a_k \in [0, 1]$  はステップ幅)

$$\hat{Q}^\pi(x_k, u_k) \leftarrow \hat{Q}^\pi(x_k, u_k) + a_k \left( l_k + \beta \hat{Q}^\pi(x_{k+1}, \pi(x_{k+1})) - \hat{Q}^\pi(x_k, u_k) \right)$$

- ベルマン方程式  $Q^*(x, u) = (\ell + \beta \mathcal{A}_x[V^*])(x, u)$  の解  $Q^*$  を求める場合は、

$$\hat{Q}^*(x_k, u_k) \leftarrow \hat{Q}^*(x_k, u_k) + a_k \left( l_k + \beta \inf_{u \in U(x_{k+1})} \hat{Q}^*(x_{k+1}, u) - \hat{Q}^*(x_k, u_k) \right)$$

- このように推定値をもとに別の推定値を構築することを**ブートストラップ** (bootstrap) とよぶ。

## アルゴリズム例 (4) SARSA と Q 学習

- **SARSA 法** :  $Q^\pi$  を推定し, 制御則は  $\pi(x) = \arg \min_{u \in U(x)} \hat{Q}^\pi(x, u)$ . 方策オン型.
- **Q 学習 (Q learning)** :  $Q^*$  を推定し, 制御則は  $\pi^*(x) = \arg \min_{u \in U(x)} \hat{Q}^*(x, u)$ . 方策オフ型.
- Q 学習の方が楽観的.

---

### アルゴリズム 9 SARSA と Q 学習

---

**input**  $\tilde{x}, \tilde{u}, l, x, a, \hat{Q}$  ▷  $(x_{k-1}, u_{k-1}, l_{k-1}, x_k)$ , ステップ幅  $a_k > 0$   
**output**  $x, u, \hat{Q}$  ▷  $x_k, u_k$   
 1: **if** type = SARSA **then**  
 2:      $u^*$  を確率  $(1 - \varepsilon)$  で  $\arg \min_{u \in U(x)} \hat{Q}(x, u)$ , 確率  $\varepsilon$  で一様乱数 ▷  $u^* \sim \pi(x_k)$   
 3: **else if** type = Q-Learning **then**  
 4:      $u^* \leftarrow \arg \min_{u \in U(x)} \hat{Q}(x, u)$  ▷  $u^* := \arg \min_{u \in U(x_k)} \hat{Q}(x_k, u)$   
 5:  $\hat{Q}(\tilde{x}, \tilde{u}) \leftarrow \hat{Q}(\tilde{x}, \tilde{u}) + a(l + \beta \hat{Q}(x, u^*) - \hat{Q}(\tilde{x}, \tilde{u}))$  ▷ 式 (9.36) または式 (9.37)  
 6:  $u$  を確率  $(1 - \varepsilon)$  で  $\arg \min_{u \in U(x)} \hat{Q}(x, u)$ , 確率  $\varepsilon$  で一様乱数 ▷  $\varepsilon$ -greedy に制御入力を計算

---

## アルゴリズム例 (5) TD 学習にもとづく方策勾配法

### actor-critic 法

1. 制御則  $\pi_\theta$  に従って,  $\bar{k}$  ステップの制御を  $\bar{s}$  巡行った標本  $(x_k^s, u_k^s)_{k=0, \dots, \bar{k}}^{s=1, \dots, \bar{s}}$  を取得
2. 制御則評価の役割を担う行動価値関数 (critic ともよばれる) をもとに,

$$\nabla_\theta J(\pi_\theta) \simeq \frac{1}{\bar{s}} \sum_{s=1}^{\bar{s}} \sum_{k=0}^{\bar{k}} \beta^k (\nabla_\theta \log(\pi_\theta(u_k^s | x_k^s))) \left( \hat{Q}^{\pi_\theta}(x_k^s, u_k^s) - b(x_k^s) \right) \quad (9.38)$$

から標本近似で勾配を計算

3. 勾配法により制御則 (actor ともよばれる) を改善

- 制御則  $\pi_\theta$  は設計者が決定できる関数であるため, その勾配も既知
- **ベースライン**  $b$  は期待値には影響せず, 分散を減らすことができ,  $b$  を状態価値関数とした  $Q^\pi - V^\pi$  はアドバンテージ関数とよばれる.

## 第 9 章

### 1. 最適制御と強化学習

### 2. モデルベース設計

ベルマン方程式

方策勾配法

### 3. モデルフリー設計

行動価値関数の推定

アルゴリズム例