

part 1

---

---

---

---

---

---



# 物理屋のための機械学習講義 8

## 。強化学習

0. 強化学習の問題設定  $s \in S \quad a \in A$

{ 環境 :  $\mathcal{E} = \{S, A, \underline{r}(\{s_t\}, \{a_t\}), \underline{p}_T(\{s_{t+1}\} | \{s_t\}, \{a_t\})\}$   
エージェント :  $ag = \{\underline{\pi}(\{a_t\} | \{s_t\}, \{a_{t-1}\})\}$

$S_0 \quad \underline{R} = \sum_{t=0}^{\infty} \gamma^t r(\{s_t\}, \{a_t\}) \quad \{a_t\} = \{a_1, \dots, a_t\}$   
を最大化するような  $\pi$  を見つけよう...

$\max_{\pi} \mathbb{E}_{\pi} [R | S_0] \dots$  目的関数

{ エージェントが  $r(\{s_t\}, \{a_t\})$  と  $p_T(\{s_{t+1}\} | \{s_t\}, \{a_t\})$  を知ってるかどうか?  
→ 知ってる : プランニング問題  
→ 知らない : (狭義の) 強化学習

---

$\begin{cases} \underline{r}(\{s_t\}, \{a_t\}) = r(s_t, a_t) \\ \underline{p}_T(\{a_{t+1}\} | \{s_t\}, \{a_t\}) = p_T(a_{t+1} | s_t, a_t) \end{cases} \Rightarrow$  Markov 的  
環境

$\pi(a_t | \{s_t\}, \{a_{t-1}\}) = \pi(a_t | s_t)$  Markov 的方針

以降は簡単のため、環境が Markov の場合に話を絞る。

# 0.1 Markov 方策の十分性

\* 環境が Markov かつ 目的関数が各時刻に同じ確率で書けるならば、Markov 方策を考えた方がよい。

step 1

$$\Pr(S_t = s, A_t = a \mid s_0, \pi^H) = \Pr(S_t = s, A_t = a \mid s_0, \pi^M) \quad \text{--- ①}$$

$$\pi_t^M := \frac{\Pr(S_t = s, A_t = a \mid s_0, \pi^H)}{\Pr(S_t = s \mid s_0, \pi^H)} \quad \text{--- ②}$$

を  $t=1, 2, \dots$   
(帰納法を使う)

$t=0$  のとき  $\Pr(S_0 = s_0 \mid s_0, \pi^H) = 1 = \Pr(S_0 = s_0 \mid s_0, \pi^M)$

① を満たすのは、  
この式が成り立つ

$t=k$  のとき ① が成り立つと仮定

$t=k+1$  のとき

$$\Pr(S_{k+1} = s, A_{k+1} = a \mid s_0, \pi^H) = \underbrace{\pi_{k+1}^M(a|s)}_{\pi^H} \Pr(S_{k+1} = s \mid s_0, \pi^H)$$

$$\begin{aligned} &= \sum_{s', a'} \underbrace{P_T(s|a, s')} \cdot \Pr(S_k = s', A_k = a' \mid s_0, \pi^H) \\ &= \sum_{s', a'} P_T(s|a, s') \times \Pr(S_k = s', A_k = a' \mid s_0, \pi^M) \quad (\text{全て Markov である}) \\ &:= \Pr(S_{k+1} = s, A_{k+1} = a \mid s_0, \pi^M) \end{aligned}$$

$$\underline{f(\pi, s_0)} = \underline{f(\Pr(S_t, A_t \mid \pi^H, s_0))}$$

ここで書ける

ex)  $\mathbb{E}_\pi [R \mid s_0]$

Markov 方策  $\pi^H$  で  $f$  を最大化可能

# マルコフ決定過程問題

$$R = \sum_{t=0}^{\infty} r(s_t, a_t)$$

$$\mathbb{E}_{\pi} [R | s_0]$$

定常にはこれを最大化したいが。  
代わりに以下を考へよ。

$$C = \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \quad (0 < \gamma < 1)$$

(本当は  $\gamma = 1$  を考えたい  
が、 $\gamma < 1$  の方が性質が  
良い)

$$V^{\pi}(s_0) = \mathbb{E}_{\pi} [C | s_0] \quad \dots \text{価値関数}$$

$$V^*(s_0) = \max_{\pi} \mathbb{E}_{\pi} [C | s_0] \quad \dots \text{最適価値関数}$$

$$= \max_{\pi_0} \mathbb{E}_{\pi_0} \left[ r(s_0, a_0) + \gamma \sum_{s_1} p_{\tau}(s_1 | s_0, a_0) V^*(s_1) \right]$$

"  $X(s, a)$

$$= \max_{a_0} \left[ r(s_0, a_0) + \gamma \sum_{s_1} p_{\tau}(s_1 | s_0, a_0) V^*(s_1) \right]$$

$$= B^* V^*(s_0)$$

$$B^* v(s_0) = \max_{a_0} \left[ r(s_0, a_0) + \gamma \sum_{s_1} p_{\tau}(s_1 | s_0, a_0) v(s_1) \right]$$

... 最適ベルマン作用素

$$\pi^*(a | s) = \delta_{a, a^*(s)} \quad a^*(s) = \underset{a}{\operatorname{argmax}} X(s, a)$$

$\Rightarrow V^*$  がわかれば、 $\pi^*$  (目的関数を最大化する最適方針) がわかる。

$\Rightarrow V^*$  を考へる上では、 $\max_{\pi} \mathbb{E}_{\pi}$  を  $\max_a$  に置きかちよ。

定常的な Markov 方針を考へれば十分最大化できる!

(前節では、Markov 方針の十分性を示した。方針の時間依存性を残す。)

◦ ベルマン作用素の縮小性

$$\underline{u}(s) - V^\pi(s) = d(s)$$

$$\underline{V}^\pi = B_\pi \underline{V}^\pi$$

[  $\pi$  ]

$$\begin{aligned} u'(s) &= \underline{B}_\pi u = \sum_a \pi(a|s) \{ r(s,a) + \gamma \sum_{s'} p_\tau(s'|s,a) \underline{V}^\pi + d \} \\ &= \underline{V}^\pi + \gamma \sum_a \pi(a|s) p_\tau(s'|s,a) \underline{d}(s') \end{aligned}$$

$$u'(s) - V^\pi(s) = \gamma \sum_a \pi(a|s) p_\tau(s'|s,a) \underline{d}(s')$$

$$\begin{aligned} |u' - V^\pi|(s) &= \left| \gamma \sum_a \pi(a|s) p_\tau(\dots) \underline{d}(s') \right| \\ &\leq \gamma \sum_a \pi(a|s) p_\tau(s'|s,a) |d(s')| \\ &\leq \underbrace{\gamma \sum_{a,s'} \pi(a|s) p_\tau(\dots)}_{=1} \max_s |d(s)| \end{aligned}$$

$$\max_s |u'(s) - V^\pi(s)| \leq \gamma \max_s |d(s)|$$

$$\underline{u} - \underline{V}^\pi = \underline{d}(s)$$

$$u' = B_\pi u$$

$$|u' - V^\pi| \leq |d(s)| \Rightarrow \gamma \max |d(s)|$$

$\pi \rightarrow \pi^*$   $B^* \dots$  同様

⇒ 縮小性から解の唯一性も言える。 (動的計画法の問題の場合)

⇒ 縮小性があるのて、ベルマン作用素が分かれば、価値関数も計算可

◦ ベルマン方程式と最小作用の原理

$$V^* = \max_a \left[ r(s, a) + \gamma \sum_{s'} p_{\tau}(s' | s, a) V^*(s') \right]$$

$$\gamma = 1. \quad p_{\tau}(s' | s, a) = \delta_{s', F(s, a)}$$

$$\underline{s'} = F(s, a) = \underline{s} + \underline{f(a)} \delta t \quad S(x, a)$$

$$s = \underline{x}$$

$$x' - x = f(a) \delta t$$

$$r(s, a) = \underline{L(x, a)}$$

$$f(a) = \dot{x}$$

$$\frac{\partial S(x, a)}{\partial a} = 0 \Rightarrow \frac{\partial L}{\partial a} + \frac{\partial}{\partial a} \left\{ V^*(x') \right\} = 0$$

$$\frac{\partial V^*}{\partial x'} \frac{\partial f}{\partial a} \delta t$$

$$\underline{V^*} = \max_{\pi} \mathbb{E}_{\pi} \left[ \sum_{\tau} r_{\tau} | x \right]$$

$$r(x, \dot{x}) = \frac{1}{2} m \dot{x}^2 - \frac{1}{2} k x^2$$

$$\int_{\tau} \frac{\partial L(x_{\tau}, a_{\tau})}{\partial a} \frac{\partial f}{\partial a} \delta t$$

$$\int dt \frac{\partial L}{\partial x} \frac{\partial f}{\partial a}$$

$$a = \dot{x}$$

$$\Rightarrow \frac{\partial L}{\partial \dot{x}} + \int dt \frac{\partial L}{\partial x} = 0 \Rightarrow \frac{d}{dt} \left( \frac{\partial L}{\partial \dot{x}} \right) + \frac{\partial L}{\partial x} = 0$$