

## マルコフゲームへの学習アルゴリズムの応用

新大 理学部 田中謙輔

### §1. マルコフゲーム

ここでは離散時間の二人零和マルコフゲームを次のような一組の組  $(S, A, B, \mu, r, \beta)$  で定義する: (1)  $S = \{1, 2, \dots, s\}$  はシステムの状態空間と呼ぶ有限集合, (2)  $A = \{a^1, a^2, \dots, a^m\}$  はプレイヤー I が選択できる純戦略の集りでプレイヤー I の行動空間と呼ぶ有限集合, (3)  $B = \{b^1, b^2, \dots, b^{m_2}\}$  はプレイヤー II が選択できる純戦略の集りでプレイヤー II の行動空間と呼ぶ有限集合, (4)  $\mu(l'|l, a, b)$  は  $l \in S, a \in A, b \in B$  によって定まる, 次の段階で状態  $l'$  がおこる推移確率, (5)  $r(l, a, b)$  は  $S \times A \times B$  の上で定義されているプレイヤー I の利得関数と呼ぶ有限な実数値関数, (6)  $\beta$  は割引因子と呼ぶ  $0 \leq \beta < 1$  なる実数.

このマルコフゲームでは二人のプレイヤーが最初の段階でシステムの状態  $l_0 \in S$  を観測し, 互に独立に戦略  $a_0 \in A, b_0 \in B$

を選択する。この結果としてプレイヤー I は利得  $r(l_0, a_0, b_0)$ , プレイヤー II は利得  $-r(l_0, a_0, b_0)$  を受ける。次にシステムの状態  $l_0$  は推移確率  $P(l_1 | l_0, a_0, b_0)$  にしたがって新しい状態  $l_1$  に移る。次に 1 段階で 2 人のプレイヤーはシステムの状態  $l_1$  を観測し、互に独立に戦略  $a_1, b_1$  を選択する。この結果としてプレイヤー I は利得  $\beta r(l_1, a_1, b_1)$ , プレイヤー II は利得  $-\beta r(l_1, a_1, b_1)$  を受ける。次にシステムの状態  $l_1$  は推移確率  $P(l_2 | l_1, a_1, b_1)$  にしたがって新しい状態  $l_2$  に移る。以下同様にして 2 人のプレイヤーはゲームを無限に続けることにする。このとき各プレイヤーが政策  $\pi, \rho$  を用いたとき、プレイヤー I の全期待割引利得は

$$V(l_0, \pi, \rho) = E^{\pi, \rho} \left[ \sum_{n=0}^{\infty} \beta^n r(l_n, a_n, b_n) \right]$$

となる。ただし  $\pi = (f_0, f_1, \dots)$ ,  $\rho = (g_0, g_1, \dots)$  で  $f_n, g_n$  は  $n$  段階で各プレイヤーが戦略を選択する写像とする。このときプレイヤー I は  $V(l_0, \pi, \rho)$  をできるだけ大きくするように、プレイヤー II はできるだけ小さくするように互に政策  $\pi, \rho$  を選択する基準で考えることにする。

このようなマルコフゲームに対して、ベクトル空間からベクトル空間への縮小写像に関する不動点定理から式:

$$\begin{aligned} V_l^* &= \max_P \min_Q \left\{ r(l, P, Q) + \beta \sum_{l'=1}^S V_{l'}^* P(l' | l, P, Q) \right\} \quad (1.1) \\ &= r(l, P^*(l), Q^*(l)) + \beta \sum_{l'=1}^S V_{l'}^* P(l' | l, P^*(l), Q^*(l)), \quad l=1, 2, \dots, S, \end{aligned}$$

$$\begin{aligned} \text{ただし } r(l, P, q) &= \sum_{i=1}^{M_1} \sum_{j=1}^{M_2} r(l, a^i, b^j) p_i q_j, \\ \phi(l|l, P, q) &= \sum_{i=1}^{M_1} \sum_{j=1}^{M_2} \phi(l'|l, a^i, b^j) p_i q_j, \end{aligned}$$

が成立する。このとき  $V^* = (V_1^*, V_2^*, \dots, V_S^*)$  はマルコフゲームのゲーム値で,  $P^*(l), q^*(l)$  がそれぞれシステムの状態  $l$  におけるプレイヤー I, プレイヤー II の最適混合戦略となることが知られている。よって二人のプレイヤーの最適政策は定常政策  $\pi^* = (P^*, P^*, \dots), \sigma^* = (q^*, q^*, \dots)$  となることがわかる。さらにゲーム値  $V^*$  は縮小写像の原理より唯一存在し次のように逐次的に求めることができる:

$$\begin{aligned} V^{(0)}: & \text{任意に与えられる } S \text{ 次元ベクトル} \\ V_l^{(n+1)} &= \max_P \min_q \left\{ r(l, P, q) + \beta \sum_{l'=1}^S V_{l'}^{(n)} \phi(l'|l, P, q) \right\} \\ & \quad l=1, 2, \dots, S. \end{aligned}$$

このとき不等式:

$$\| V^{(n)} - V^* \| \leq \beta^n \| V^{(0)} - V^* \| \quad (1.2)$$

が成立する。ただし  $V^{(n)} = (V_1^{(n)}, V_2^{(n)}, \dots, V_S^{(n)})$  で,  $\beta$  は各要素の絶対値の最大値を表している。

## §2. dummy game の正規化

マルコフゲームについての説明から dummy game と呼ばれる, システムの各状態  $l$  で利得行列  $\{V_l^*(a, b); a \in A, b \in B\}$  をもつ二人零和ゲームの最適混合戦略の組, すなわち, 二人のプ

プレイヤーの任意の混合戦略  $P, q$  に対して

$$V_{\ell}^*(P^*(\ell), q) \geq V_{\ell}^* \geq V_{\ell}^*(P, q^*(\ell))$$

が成立する  $(P^*(\ell), q^*(\ell))$  を求めることになる。ただし

$$V_{\ell}^*(a, b) = r(\ell, a, b) + \beta \sum_{\ell'=1}^S V_{\ell'}^* P(\ell' | \ell, a, b),$$

$$V_{\ell}^*(P, q) = \sum_{i=1}^{m_1} \sum_{j=1}^{m_2} V_{\ell}^*(a^i, b^j) p_i q_j.$$

ここではこのようなゲームを正規化するために利得行列を、 $\delta > 0$  に対して、 $\{V_{\ell}^*(a, b) - \frac{\delta}{2}(p_a - q_b); a \in A, b \in B\}$  と修正する。

このときプレイヤー I の期待利得は

$$V_{\ell, \delta}^*(P, q) = V_{\ell}^*(P, q) - \frac{\delta}{2}(\|P\|^2 - \|q\|^2)$$

で与えられる。ただし  $P = (p_1, p_2, \dots, p_{m_1})$ ,  $q = (q_1, q_2, \dots, q_{m_2})$  で、

$\|\cdot\|$  はユークリッドノルムを表している。したがってプレイヤー II の期待利得は  $-V_{\ell, \delta}^*(P, q)$  で与えられる。

次に、各プレイヤーが利用できる混合戦略  $P, q$  はそれぞれ  $S_{\varepsilon}^{m_1}, S_{\varepsilon}^{m_2}$  の中にあると仮定する。ただし

$$S_{\varepsilon}^m = \left\{ x = (x_1, x_2, \dots, x_m); x_i \geq \varepsilon, (i=1, 2, \dots, m), \sum_{i=1}^m x_i = 1, (0 \leq \varepsilon \leq \frac{1}{m}) \right\}.$$

このように制限された二人零和ゲームでは、 $\delta > 0$  に対して

$V_{\ell, \delta}^*(P, q)$  は厳密に凸となる。よって任意の  $P \in S_{\varepsilon}^{m_1}, q \in S_{\varepsilon}^{m_2}$  に対して

$$V_{\ell, \delta}^*(P^*(\ell, \varepsilon, \delta), q) \geq V_{\ell, \delta}^*(P^*(\ell, \varepsilon, \delta), q^*(\ell, \varepsilon, \delta)) \geq V_{\ell, \delta}^*(P, q^*(\ell, \varepsilon, \delta)) \quad (2.1)$$

となる唯一の最適混合戦略の組  $(P^*(\ell, \varepsilon, \delta), q^*(\ell, \varepsilon, \delta))$  が存在する。

このような最適混合戦略の粗について次のような2つの補助定理が重要となる。

補助定理1.

数列  $\{\varepsilon(n)\}, \{\delta(n)\}$  が次のような条件をみたしているとする:

$$\varepsilon(n) \in (0, \hat{\varepsilon}), \quad \delta(n) > 0, \quad \lim_{n \rightarrow \infty} \varepsilon(n) = \lim_{n \rightarrow \infty} \delta(n) = 0,$$

$$\lim_{n \rightarrow \infty} \frac{\varepsilon(n)}{\delta(n)} = \mu \in [0, \infty).$$

このとき  $(P^*(l, \varepsilon(n), \delta(n)), q^*(l, \varepsilon(n), \delta(n)))$  は最初のゲームの鞍部点  $(P^*(l), q^*(l))$  に収束する。

補助定理2.

任意の利得行列に対して、次のような条件をみたす  $\delta \in (0, \infty)$  と定数  $K_1, K_2, K_3$  が存在する:  $\forall \varepsilon_1, \varepsilon_2 \in [0, \hat{\varepsilon}], \forall \delta_1, \delta_2 \in (0, \delta)$ ,

$$\begin{aligned} & \|P^*(l, \varepsilon_1, \delta_1) - P^*(l, \varepsilon_2, \delta_2)\| + \|q^*(l, \varepsilon_1, \delta_1) - q^*(l, \varepsilon_2, \delta_2)\| \\ & \leq K_1 |\varepsilon_1 - \varepsilon_2| + K_2 |\delta_1 - \delta_2| + K_3 \left| \frac{\varepsilon_1}{\delta_1} - \frac{\varepsilon_2}{\delta_2} \right|. \quad (2.2) \end{aligned}$$

### §3. 学習アルゴリズム

ここでは各プレイヤー(生徒)がマルコフゲームにおける利得関数, システムの状態を変化させる推移確率を知っていないとする。このとき dummy game を通してホウ着(教師)が各プレイヤーに最適戦略  $P^*(l), q^*(l)$  をそれぞれ学習させる方法を次のように考へることにする。

まずホウ着は初期ベクトル  $V^{(0)} = (V_1^{(0)}, V_2^{(0)}, \dots, V_s^{(0)})$  を定め,

各プレイヤーは次の段階の各システムの状態で用いる初期混合戦略  $P^{(0)}(l), q^{(0)}(l), l=1, 2, 3, \dots, S$ , を定める。今  $n$  段階で到着はゲーム値の近似ベクトル  $V^{(n)} = (V_1^{(n)}, V_2^{(n)}, \dots, V_S^{(n)})$  を定めており、各プレイヤーは次の段階の各状態  $l$  で用いる混合戦略  $P^{(n)}(l), q^{(n)}(l)$  をそれぞれ定めるとする。このとき  $(n+1)$  段階では、各状態  $l$  で各プレイヤーは  $P^{(n)}(l), q^{(n)}(l)$  でゲームをプレイする。この結果各プレイヤーが用いた純戦略をそれぞれ  $x_{n+1}(l), y_{n+1}(l)$  とするとき、到着は各プレイヤーに

$$V_l^{(n+1)}(x_{n+1}(l), y_{n+1}(l)) = T(l, x_{n+1}(l), y_{n+1}(l)) + \beta \sum_{l'=1}^S V_{l'}^{(n)} \phi(l'|l, x_{n+1}(l), y_{n+1}(l))$$

の値を教える。これを用いて各プレイヤーは次の段階で用いる混合戦略を次のように作ることにする。プレイヤー I はシステムの各状態  $l$  で

$$P^{(n+1)}(l) = \Pi_{S_{\varepsilon(n+1)}^{m_1}} [P^{(n)}(l) + \delta(n+1) A^{(l)}(x_{n+1}(l), y_{n+1}(l))] \quad (3.1)$$

と作る。ただし  $A^{(l)}(a^i, b^j)$  の  $k$  要素  $A_k^{(l)}(a^i, b^j)$  は

$$A_k^{(l)}(a^i, b^j) = \begin{cases} \frac{V_l^{(n+1)}(a^i, b^j)}{\phi_i^{(n)}(l)} - \delta(n+1), & k = a^i \\ -\frac{1}{m_1 - 1} \left( \frac{V_l^{(n+1)}(a^i, b^j)}{\phi_i^{(n)}(l)} - \delta(n+1) \right), & k \neq a^i \end{cases}$$

で、記号  $\Pi_S[X]$  は有界閉集合  $S$  への  $X$  の射影作用素を表している。

同様にして、プレイヤー II はシステムの各状態  $l$  で

$$q_j^{(n+1)}(l) = \prod_{s \in \mathcal{E}(n+1)}^{m_2} [q_j^{(n)}(l) - \delta(n+1) B^{(l)}(\alpha_{n+1}(l), y_{n+1}(l))] \quad (3.2)$$

と作る。ただし  $B^{(l)}(a^i, b^j)$  の  $k$  要素  $B_k^{(l)}(a^i, b^j)$  は

$$B_k^{(l)}(a^i, b^j) = \begin{cases} \frac{V_k^{(n+1)}(a^i, b^j)}{q_j^{(n)}(l)} + \delta(n+1) & , k = b^j \\ -\frac{1}{m_2 - 1} \left( \frac{V_k^{(n+1)}(a^i, b^j)}{q_j^{(n)}(l)} + \delta(n+1) \right) & , k \neq b^j \end{cases}$$

で表している。

また次の着は次の段階で用いるゲーム値の近似ベクトル  $V^{(n+1)} = (V_1^{(n+1)}, V_2^{(n+1)}, \dots, V_s^{(n+1)})$  を次のように作る。

$$V_k^{(n+1)} = \max_P \min_Q \left\{ r(l, P, Q) + \beta \sum_{j=1}^s V_k^{(n)} p(l^j | l, P, Q) \right\}.$$

このようにして  $(n+1)$  段階での着の作業と各プレイヤーの学習は終る。以下の各段階で同様にして各プレイヤーの学習は進行する。

このとき次の定理1, 定理2によって各プレイヤーは学習の目的を達成することが保証されている。

定理1.

数列  $\{\varepsilon(n)\}, \{\delta(n)\}, \{\gamma(n)\}$  が次のような条件をみたしている:

$$(a) \quad \delta(n) > 0, \quad \gamma(n) > 0, \quad \varepsilon(n) \in (0, \hat{\varepsilon}), \quad n=1, 2, \dots,$$

$$\delta(n) \rightarrow 0, \quad n \rightarrow \infty, \quad \text{ただし } \hat{\varepsilon} = \min\left(\frac{1}{m_1}, \frac{1}{m_2}\right),$$

$$(b) \quad \lim_{n \rightarrow \infty} \varepsilon(n)/\delta(n) = \mu < \infty, \quad (c) \quad \sum_{n=1}^{\infty} \gamma(n)\delta(n) = \infty,$$

$$\begin{aligned}
 (d) \quad & \sum_{n=1}^{\infty} \gamma^2(n) \delta^2(n) < \infty, \quad (e) \quad \sum_{n=1}^{\infty} \gamma^2(n) / \varepsilon(n-1) < \infty, \\
 (f) \quad & \sum_{n=1}^{\infty} |\varepsilon(n) - \varepsilon(n-1)| < \infty, \quad (g) \quad \sum_{n=1}^{\infty} |\delta(n) - \delta(n-1)| < \infty, \\
 (h) \quad & \sum_{n=1}^{\infty} \left| \frac{\varepsilon(n)}{\delta(n)} - \frac{\varepsilon(n-1)}{\delta(n-1)} \right| < \infty.
 \end{aligned}$$

このとき, 任意の初期混合戦略の組  $(P^{(0)}(l), q^{(0)}(l)) \in S_{\varepsilon^{(0)}}^{m_1} \times S_{\delta^{(0)}}^{m_2}$  に対して, 列  $(P^{(n)}(l), q^{(n)}(l))$  は  $n \rightarrow \infty$  のとき  $(P^*(l), q^*(l))$  に確率1で収束する。

証明. 補助定理2を用いて,  $\forall n \geq n_0$  に対して

$$\begin{aligned}
 E[d(n+1) | P^{(n)}(l), q^{(n)}(l), l=1,2,\dots,s] &\leq (1 - L_1 \gamma(n+1) \delta(n+1)) d(n) + \\
 &+ K_1 |\varepsilon(n+1) - \varepsilon(n)| + K_2 |\delta(n+1) - \delta(n)| + K_3 \left| \frac{\varepsilon(n+1)}{\delta(n+1)} - \frac{\varepsilon(n)}{\delta(n)} \right| + \\
 &+ K_4 \beta^{n+1} + K_5 \frac{\gamma^2(n+1)}{\varepsilon(n)} + K_6 \gamma^2(n+1) \delta^2(n+1)
 \end{aligned}$$

が成立するように  $n_0$  と定数  $K_1, K_2, K_3, K_4, K_5, K_6$  が存在する。

以下" L

$$d(n+1) = \sum_{l=1}^s \left( \frac{m_1-1}{m_1} \|P^{(n+1)}(l) - P^{(n+1)*}(l)\|^2 + \frac{m_2-1}{m_2} \|q^{(n+1)}(l) - q^{(n+1)*}(l)\|^2 \right).$$

ここで  $D(n) = d(n) + \sum_{k=n+1}^{\infty} \beta(k)$  とおく。以下" L

$$\begin{aligned}
 \beta(n+1) &= K_1 |\varepsilon(n+1) - \varepsilon(n)| + K_2 |\delta(n+1) - \delta(n)| + K_3 \left| \frac{\varepsilon(n+1)}{\delta(n+1)} - \frac{\varepsilon(n)}{\delta(n)} \right| + \\
 &+ K_4 \beta^{n+1} + K_5 \gamma^2(n+1) / \varepsilon(n) + K_6 \gamma^2(n+1) \delta^2(n+1).
 \end{aligned}$$

このとき  $\{D(n)\}$  は semi-martingale と取り, semi-martingale に関する収束定理と  $\sum_{n=0}^{\infty} \gamma(n+1) \delta(n+1) E[d(n)] < \infty$  が成立することよりこの定理の証明は完成できる。



次に二乗平均収束を示すために次の補助定理が重要である。  
補助定理3.

非負数列  $\{u(n)\}$  が次のような条件をみたしているとする:

$$\forall n \geq n_0$$

$$(a) \quad u(n) \leq u(n-1)(1-\lambda(n)) + \theta(n),$$

$$(b) \quad \lambda(n) \in (0, 1], \quad \sum_{n=1}^{\infty} \lambda(n) = \infty, \quad \lim_{n \rightarrow \infty} \frac{\theta(n)}{\lambda(n)} = 0.$$

このとき  $\lim_{n \rightarrow \infty} u(n) = 0$  が成立する。

定理2.

定理1における条件 (d) ~ (h) が次のような条件に置き換え  
られている:

$$(d') \quad \gamma(n+1)\delta(n+1) \rightarrow 0 \quad \text{as } n \rightarrow \infty,$$

$$(e') \quad \frac{\gamma(n+1)}{\varepsilon(n)\delta(n+1)} \rightarrow 0 \quad \text{as } n \rightarrow \infty,$$

$$(f') \quad \frac{1}{\gamma(n+1)\delta(n+1)} |\varepsilon(n+1) - \varepsilon(n)| \rightarrow 0 \quad \text{as } n \rightarrow \infty,$$

$$(g') \quad \frac{1}{\gamma(n+1)\delta(n+1)} |\delta(n+1) - \delta(n)| \rightarrow 0 \quad \text{as } n \rightarrow \infty,$$

$$(h') \quad \frac{1}{\gamma(n+1)\delta(n+1)} \left| \frac{\varepsilon(n+1)}{\delta(n+1)} - \frac{\varepsilon(n)}{\delta(n)} \right| \rightarrow 0 \quad \text{as } n \rightarrow \infty.$$

さらに、次の条件がみたされている:

$$(i) \quad \frac{\gamma^{n+1}}{\gamma(n+1)\delta(n+1)} \rightarrow 0 \quad \text{as } n \rightarrow \infty.$$

このとき  $(P^{(n+1)}(l), q^{(n+1)}(l))$  は二乗平均の意味で  $(p^*(l), q^*(l))$  に収束する。

(注意) アルゴリズムにおける数列の例として

$$\gamma(n) \sim \frac{1}{n^\alpha}, \quad \varepsilon(n) \sim \frac{1}{n^\nu}, \quad \delta(n) \sim \frac{1}{n^\sigma},$$

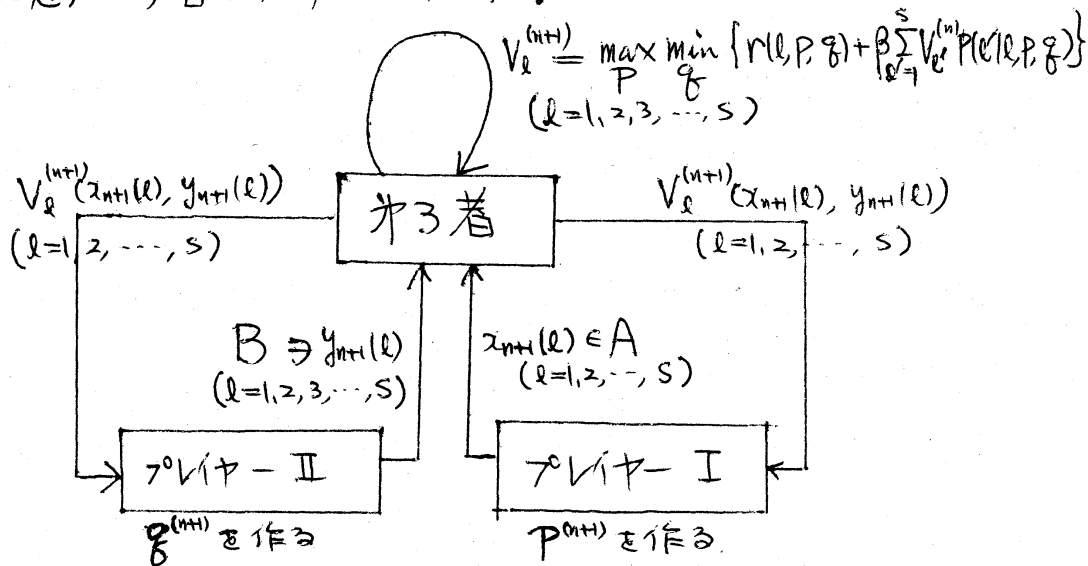
$$\frac{\varepsilon(n)}{\delta(n)} - \mu \sim \begin{cases} \frac{1}{n^\nu} & , \beta = \sigma \\ \frac{1}{n^{\beta-\sigma}} & , \beta > \sigma, \end{cases}$$

とおき定理1に対しては,  $\beta \geq \sigma > 0$ ,  $\nu > 0$ ,  $\frac{1}{2} < \alpha + \sigma \leq 1$ ,

$2\alpha - \beta > 1$ , 定理2に対しては  $\beta \geq \sigma > 0$ ,  $\nu > 0$ ,  $\alpha + \sigma \leq 1$ ,

$\alpha - \beta - \sigma > 0$  等に選択すればよい。

(注意) 学習システムの図示。



### 参考文献

1. Tsypkin, Ya. Z., *Adaptation and Learning in automatic systems*, (1968) (in Russian).
2. Nagin, A. Z. and Poznyak, A. S., *Stochastic zero-sum game of two automata*, *Avtom. Telemekh* (1977) (in Russian).
3. Tanaka, K. and Homma, H., *On the learning algorithm of 2-person zero-sum game*, to appear.