

学習プロセスと後悔

大阿久 博

(武蔵野大学政治経済学部教授)

一 学習理論

従来のゲーム理論において中心となってきた解概念はナッシュ均衡であり、ナッシュ均衡が一般的に存在することはナッシュ自身によって証明された (Nash[89])。しかし、均衡が(一般には複数)存在すること、それらを具体的に求めること、さらにそれらの中のどの均衡が実現するか見極めることは、決して同じことではない。様々な実験を通して、こうしたことが明らかにされてきた。つまりゲーム理論においては、「均衡が具体的に求められない」という非決定性の問題」と「均衡を計算して求めるのは原理的に可能であるが、実際には膨大な時間がかかって困難である」という非決定性の問題」という異なる問題があるのである。(川越[12]第三章、p.79)。後者の問題が生じたとき、人間(プレイヤー)はどのように均衡に達するのかを研究するのが学習理論である。

こうした学習理論においては、プレイヤーは従来のゲーム理論が想定してきたような合理性は持つておらず、必然的に「限定合理的」になる。¹⁾

学習理論では、主に二つのタイプの学習プロセスが考えられている。一つは「信念学習」、もう一つは「強化学習」と呼ばれるものである。前者は、相手プレイヤーが過去に選択した戦略の頻度から、次に相手がとる戦略を予想し、その予想に対して最適な戦略を選択するという学習行動であり、後者は、相手のプレイヤーの戦略の選択頻度ではなく、過去に選んだ自分自身の戦略がどの程度の利得を上げてきたかを考え、より利得の高かったものをより高い頻度で選択するという学習行動である。信念学習はゲーム理論においても「仮想プレイ (Fictitious play)」などで議論されてきた。一方、強化学習は心理学などで研究されてきた分野である。

以下では、「強化学習」の考え方、その中でも特に「後悔しない」戦略を選択するといった行動様式 (regret matching モデル) について既存研究を概観する。²⁾

1.1 後悔 (regret)

1. 例

「後悔しない戦略をとる」という状況を定式化するに当たり、例として次の表1のような 2×2 ゲームを考へよう (The Soda Game (Young[3])). このゲームでは、二人のプレイヤーはそれぞれ L と R の二つの行動を持ち、利得の代わりに財 (Coke, Sprite, Seven-up, Pepsi) が表示されている。

表1 The Soda Game

		Player2	
		L	R
Player1	L	Coke, Coke	Sprite, Seven-up
	R	Seven-up, Sprite	Pepsi, Pepsi

表2 Soda Game が11回繰り返されたときのプレイヤー1の利得

	1期	2期	3期	4期	5期	6期	7期	8期	9期	10期	11期	12期
Player1	L	R	R	R	L	L	L	L	R	R	R	?
Player2	L	L	R	L	R	L	R	L	R	L	L	?
1の利得	1	0	1	0	0	1	0	1	1	0	0	?

表3

	L	R
L	1, 0	0, 1
R	0, 1	1, 0

表4

	L	R
L	1, 1	0, 0
R	0, 0	1, 1

表5 Soda Game : 常にRを選択

	1期	2期	3期	4期	5期	6期	7期	8期	9期	10期	11期	12期
Player1	R	R	R	R	R	R	R	R	R	R	R	?
Player2	L	L	R	L	R	L	R	L	R	L	L	?
1の利得	0	0	1	0	1	0	1	0	1	0	0	?

表6 Soda Game : 常にLを選択

	1期	2期	3期	4期	5期	6期	7期	8期	9期	10期	11期	12期
Player1	L	L	L	L	L	L	L	L	L	L	L	?
Player2	L	L	R	L	R	L	R	L	R	L	L	?
1の利得	1	1	0	1	0	1	0	1	0	1	1	?

各々のプレイヤーはこれらの財から利得を得るが、自分自身の利得についてはわかるが、相手が得る利得についてはわからないものとする。このような相手の得る利得がわからない（利得の分布さえ分らないかもしれない）ゲームを**不確実性ゲーム**（uncertain game）と呼ぶことにする。

以下、このゲームを繰り返し行い、十一回の結果が表2のようになったとしよう。一番下の数字がプレイヤー1の各期の利得である。割引を考えなければ、プレイヤー1の利得は十一期合計で5になっている。

プレイヤー1の立場で考えたとき、十二回目にはLとRのどちらを選んだらよいか。このゲームでは相手の利得がわからないので、そもそも自分（プレイヤー1）がどんなゲームに参加しているのかわからない。協調ゲーム（表3）であるかもしれないし、コイン合わせ（matching pennies）ゲーム（表4）であるかもしれない（もちろんそれら以外かもしれない）。

ここで、プレイヤー1がすべての期で同じ行動を選択しなければならないことを条件として、このゲームをやり直せる機会が与えられたとしよう。むしろやり直しが可能になったとき、相手のプレイヤー2がどのような対応してくるかはわからないが、ここでは非常に単純に、相手は行動を変更してこないと考えよう。

このとき、プレイヤー1が「常にRを選択」（このような行動様式を R^* と記す）に変更すると、利得は表5のようになる。このときにはプレイヤー1の利得は4に減少する。

また、プレイヤー1が「常にLを選択」（ L^* とする）に変更すると、プレイヤー1の利得は7に増加する（表6）。

したがって、

- ・ R_* を取らなかったことについては後悔しないが、
- ・ J_* を取らなかったことには後悔する

であろう。よって後悔を減らすため、将来は L の選択にもっと大きなウェイトを置くべきと考えるだろう。

以下ではより一般的に、一人のプレイヤーが「自然 (Nature)」を相手にプレーする場合を考える。ここでは自然とは、相手のプレイヤーと考えてもよいし、自分が置かれている状態と考えてもよいが、いずれにせよ上述のように、意思決定を行うプレイヤーとの interactive な状況は考えない。自然の選ぶ行動 (あるいは状態) は観察できるとする。もちろん意思決定プレイヤーの利得は、自身の意思決定と自然の選ぶ行動に依存する。

2. Hannan, Fudenberg and Levine, Hart and Mas-Colell モデル

有限個の要素からなる状態の集合を Ω とする。「自然」はゲームの行われる各期において Ω から一つの状態を選ぶ。また、この自然とゲームを行う一人のプレイヤーを想定し、彼女/彼は行動の集合 A から一つを選ぶとする。このゲームが繰り返し行われ、各 τ 期において、自然が状態 $\omega_\tau \in \Omega$ を選び、プレイヤーの行動が $a_\tau \in A$ であったとき、プレイヤーは利得 $u(\omega_\tau, a_\tau)$ を得る。また、 $h_\tau = (\omega_1, \dots, \omega_\tau, a_1, \dots, a_\tau)$ で長さ τ の履歴を表す。便宜上、履歴がない初期状態を空履歴 h_0 と表し、この空履歴と有限個の履歴すべてからなる集合を H と表記する。

行動の集合 A 上の確率分布を $\Delta(A)$ とする。プレイヤーの取る戦略は H から $\Delta(A)$ への関数 $f: H \rightarrow \Delta(A)$ で与えられる。

ゲームは次のように進行する。まず、第1期に関数 f が最初の行動 a_1 （これを a_1 と書く）が従う分布を決定する $(f(h_1) \cap \Delta(A))$ 。次の第2期では、実現した状態 ε_1 と a_1 に基づいて、 f は第2期の行動 a_2 が従う分布を決定する。第3期以降も同様である。つまり $t+1$ 期の行動 a_t の従う確率分布は $f(h_t) = f(\omega_1, \dots, \omega_t; a_1, \dots, a_t) \cap \Delta(A)$ である。また、無限の状態の列 $\varepsilon = (\omega_1, \omega_2, \dots, \omega_n, \dots)$ に対して、戦略 f と ε は無限の行動の列 A_{∞} 上の確率分布を導く（この確率分布を (ω, f) と書くことにする。ただし \mathbb{Z} は自然数の集合）。

さて、今まで取ってきた戦略を後悔するかどうかは、その戦略を取らずに「他の戦略」を取ったときに利得が増加するかどうかにかかっている。「他の戦略」として、まず一番単純に過去の履歴に関係なく、ある戦略 $a \in \Delta$ をずっと取り続ける戦略 p^* を考えよう。状態の列 $\varepsilon = (\omega_1, \omega_2, \dots, \omega_n)$ に対して、 t 期の行動を a_t ($t = 1, 2, \dots, n$)とすると、 t 期における a^* との利得の差は、 $u(\omega_t, a_t) - u(\omega_t, a)$ であり、第1期から n 期までの利得の差の平均は

$$L_n(f, a^*) = \frac{\sum_{t=1}^n [u(\omega_t, a_t) - u(\omega_t, a)]}{n}$$

である。ここでゲームの回数を大きくしていったとき ($n \rightarrow \infty$)、 $\liminf_n L_n(f, a^*) \geq 0$ であるならば、 p^* を取らなかったことを後悔しないだろう。Hannan[3]によって、次のことが証明されている。

定理 1 (Hannan[3])

任意の $a \in \Delta$ と任意の状態の列 $\varepsilon = (\omega_1, \omega_2, \dots, \omega_n, \dots)$ に対して、 a^* と比べ後悔しない戦略 f が存在する。

Hannan[3]の結果は、同一の行動 a がずっと取られたときという極めてナイーブな戦略と比較した場合を扱っている。この結果は、

- (1) 期の列を互いに素な部分集合に分割し、ある部分集合上では特定の行動 $a \in \Delta$ をとる戦略、
 - (2) 二つの異なる行動 a, b に対して、 a がプレーされていた場合にはそれをすべて b で置き換えるような戦略
- に対して後悔しない戦略 f を求めるという二つの方向に拡張された。(1) については Fudenberg and Levine [2]、(2) については Hart and Mas-Colell [4] によって分析されている。Fudenberg and Levine [2] では、まず期の集合 (自然数の集合 \mathbb{N} と同一視できる) の分割 B_1, B_2, \dots, B_k を考える。そして B_j については任意の行動を a で置き換え、 B_j 以外のところでは行動を変更しないという戦略を $a^* | B_j$ と記すことにする。この a^* も Fudenberg and Levine [2] の結果は次のように書ける (Lehrer [6])。

定理 2 (Fudenberg and Levine [2])

任意の期の集合の分割 B_1, B_2, \dots, B_k について、任意の状態の列 $\omega = (\omega_1, \omega_2, \dots, \omega_n, \dots)$ と任意の $a \in A, j$ (ただし $1 \leq j \leq k$) に対して、 $a^* | B_j$ と比較して後悔しない戦略 f が存在する。

Hart and Mas-Colell [4] では、任意の行動のペア (a, b) を考え、ある期 t の行動が a であった場合それを b で置き換え、 a 以外の場合には置き換えない戦略を考える。これは関数

$$g_{a,b} : H \times A \rightarrow A - |a|,$$

$$g_{a,b} = (h_{t-1}, a) = b$$

$$g_{a,b} = (h_{t-1}, c) = c : c \neq a$$

で表せる。このとき Hart and Mas-Colell [4] (Theorem A) の結果は、次のようになる。

定理 3 (Hart and Mas-Colell [4])

任意の状態の列 $\omega = (\omega_1, \omega_2, \dots, \omega_n, \dots)$ に対して、 $g_{a,b}$ と比較して後悔しない戦略 σ が存在する。

3. Lehrer モデル

さらに Lehrer [6] では、「置換スキーム」と「アクティブ関数」という二つの関数を導入することで、Fudenberg and Levine [2] と Hart and Mas-Colell [4] の分析の一部をスペシアルケースとして含む、より一般的な学習プロセスが提案されている。

「置換スキーム」は過去の履歴において採用してきたある行動を、別の行動で置き換える関数

$$g : H \times A \rightarrow A$$

である。行動の置換には様々なケースが考えられる。前述の Hart and Mas-Colell [4] は、任意の二つの行動のペアを考え、その間で置換が起きるといふ特別の場合を扱っていると見なせる。

さらに Lehrer [6] では、プレイヤーが特定の期において戦略のパフォーマンスを重視するケースを想定する。例えばプレイヤーが、休曜日や週末には高い利得が得られるようにしたいと望むような場合である。そのため、こうした「特別な期(日)」を特定化する関数、アクティブ関数を次のように定義する。

$$I : H \times A \rightarrow \{0, 1\}$$

このアクティブ関数により、 $I(h_{t-1}, a_t) = 1$ であるなら、行動 a_t の取られた第 t 期はアクティブであるという

ことになる。履歴が $h_{t-1} \in \Omega$ であるとき、 t 期までに関数 I がアクティブになった、つまり関数の値が 1 になった回数を

$$I^n(h_{t-1}, a_t) = \sum_{s=1}^t I(h_{s-1}, a_s)$$

とする。これは履歴にしたがって t 期までに何回アクティブ期があったかを表している。ここで二つの関数のペア $(g, D) : H \times A \rightarrow A \times \{0, 1\}$ により、ある期がアクティブであるかどうか、またその期にどのような行動の置換が行われるかが定められることになる。

次に戦略 γ と、置換スキーム σ によって置き換えられた戦略が、アクティブな期を通して平均的にどちらが良いか、あるいは、戦略 γ を取ることで後悔しないかを考えるよう。したがって以下ではプレイヤーは γ と (g, D) を比較すると考えることができる。

いまある期 t において行動 a_t が取られるとする。このときの利得は $u(\omega_t, a_t)$ である。これに対して、置換スキーム σ で行動を置き換えた場合、 t 期の行動は $g(h_{t-1}, a_t)$ であるから、 t 期の利得は $u(\omega_t, g(h_{t-1}, a_t))$ となる。したがってプレイヤーにとって、 t 期において σ で行動を変えた場合との利得の差は、 $u(\omega_t, a_t) - u(\omega_t, g(h_{t-1}, a_t))$ となる。この差が生まれるのは t 期がアクティブ関数 I によってアクティブになる場合である。

以上から、第 n 期までの利得の差の平均は、

$$L_n(f, g) = \frac{\sum_{t=1}^n I(h_{t-1}, a_t) [u(\omega_t, a_t) - u(\omega_t, g(h_{t-1}, a_t))]}{I^n(h_{n-1}, a_n)}$$

となる。ゲームの回数を $\rightarrow \infty$ と考えたとき、アクティブ関数がアクティブになる回数が無限になるようなほとんどすべての行動の列に対して、 $\liminf_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T (f_t(s_t) - v)$ と $\limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T (f_t(s_t) - v)$ との間が v となる場合、戦略 γ は v を平均として得られる利得より高い、つまり戦略 γ を取ったことを後悔しないということである。このことすべての (s, v) からなる集合 R を考える。Lehrer [6] よって次が示されている。

定理 4 (Lehrer [6])

任意の状態の列 $\omega = (\omega_1, \omega_2, \dots, \omega_n, \dots)$ に対して、ほとんどすべての $(s, v) \in R$ と比べ後悔しない戦略 γ が存在する。

これら一連の「後悔なし (no regret)」の定理は Blackwell [1] の approachability 定理によって証明されている。特に Lehrer [6] では、オリジナルの Blackwell 定理を無限次元に拡張した approachability 定理 (Lehrer [5]) が使われている。⁴⁾

三. 有限記憶

次に N 人による標準形ゲームを考えよう。各プレイヤーは有限個の要素から成る行動集合 A_i を持つものとする。 $A = \prod_{i \in N} A_i$ とし、一般に各プレイヤーの行動の組は $a = (a_1, a_2, \dots) \in A$ と表われ、 a_i はプレイヤー i の行動

a_{-i} はプレイヤー i 以外のプレイヤーの行動の組を示す。今まで同様 $\mathbb{N}, \mathbb{A} \rightarrow \mathbb{R}$ (実数体) をプレイヤー i の利得関数とする。行動の組が (a_i, a_{-i}) のとき、プレイヤー i の利得は $u_i(a_i, a_{-i})$ である。

このとき $\forall i \in \mathbb{N}, a_i, b_i \in \mathbb{A}$ に対し、

$$\sum_{a_{-i}} u_i(b_i, a_{-i}) q(a_i, a_{-i}) \leq \sum_{a_{-i}} u_i(a_i, a_{-i}) q(a_i, a_{-i})$$

を満たす \mathbb{A} 上の分布 $q: \mathbb{A} \rightarrow \Delta(\mathbb{A})$ を相関均衡と呼ぶ。

このゲームが各期 $t (= 1, 2, \dots)$ 毎に繰り返し行なわれる場合を考える。 t 期の行動の組を (a_t, a_{-t}) とすると Hart and Mas-Colell [4] では、上述のように行動 a_t を b_t で置き換える場合を想定していた。ここでの置換スキームは、

$$g_{a_t, b_t}^i = (a_t^i, a_{-t}^i) = b$$

$$g_{a_t, b_t}^i = (a_t^i, a_{-t}^i) = a_t^i; a_{-t}^i \neq a_t^i$$

と書ける。

また、 t 期までの (a_t, a_{-t}) 組が生じる確率

$$\mu(a_t, a_{-t}) = \frac{1}{T} \mathbb{1}_{\|k \leq t; (a_t^i, a_{-t}^i) = (a_t, a_{-t})\|}$$

を経験分布 (empirical distribution) と呼ぶ。

Hart and Mas-Colell [4] では、 $t \rightarrow \infty$ のときこの経験分布が (ほとんど確実に) 相関均衡の集合に収束することが明らかにされた。また、経験分布が相関均衡の集合へ収束することが、後悔しない戦略が存在することと同値になることも示されている。しかしながら、経験分布が実際にどの「点」に収束するのかといった収

束の細かい特徴についてはほとんど知られていない (Hart and Mas-Colell[4] Section 2)。

この Hart and Mas-Colell[4] のモデルをはじめ、いままで考えてきたモデルではプレイヤーは履歴について完全に記憶していた。つまり、過去に行われた各期でのゲームの利得を完全に知っており、それをもとに戦略を変更するかどうか検討するという行動様式を考えてきたのである。

しかし過去の自分の行動や利得をすべて記憶しているという状況は、現実的ではない場合が多いであろう。その点に注目し、限られた過去の記憶しか持たないプレイヤーを分析したのが Saran and Serrano [10] である。こうすることによって、Hart and Mas-Colell[4] よりも、経験分布についてのシャープな結果が得られる。以下それを紹介する。

Saran and Serrano [10] では、プレイヤーが戦略を考えるに当たり初期時点からすべての期の状況を考慮に入れるのではなく、直近の日期のみ考慮する、言い換えると、直近の日期の記憶しか持たない場合を想定した。プレイヤーの (n) 期から遡って、直近日期の平均利得は、

$$L_i^m(a_i^t, a_{-i}^t) = \frac{1}{m} \sum_{k=t-m+1}^t u_i(a_i^k, a_{-i}^k)$$

となる。

次に、直近日期間の行動で、 (n) 期の行動 a_{-i}^n と同じ行動を取った期をすべて b_i で置き換える場合を考えよう。ただし b_i 以外の行動を取っている場合は変更しないものとする。したがって直近日期間の置換スキームは、Hart and Mas-Colell[4] と同じものになる。この置換スキームを使った日期の平均利得は、

$$L_i^m(b_i, a_{-i}^i) = \frac{1}{m} \sum_{k=1}^m u_i(g_{i(b_i^k, a_{-i}^k)})$$

と書ける。直近の t 期間における行動 b_i を b_i に変換したときとしないときの平均利得の差が正であった場合には、 $t+1$ 期には正の確率で行動 b_i を取り、また、そのまま a_i を取り続ける確率も正であるとする。つまり、 $L(b_i, a_i^i) = L_i^m(b_i, a_{-i}^i) - L_i^m(a_i^i, a_{-i}^i) > 0$ の場合に確率 $q(L(b_i, a_i^i) > 0)$ で行動 b_i への変更が起きる。ただし、 $\sum_{b_i \in A_i} q(L(b_i, a_i^i) > 0) > 1$ と仮定する。

このような状況で、直近 t 期の行動の組の履歴を一つの状態と考える。したがって状態の集合は Σ_t^m である。確率分布 ρ が与えられたとき、こうした有限記憶によって推移するプロセスは、 Σ_t^m 上の非周期的マルコフ・プロセス $M(q)$ になる。

任意の行動の組 $(a_i, a_{-i}) \in A$ に対して、プレイヤー i の同等以上の反応 (same or better reply) からなる集合を

$$R_i(a_i, a_{-i}) = \{b_i \in A_i \mid b_i = a_i \text{ \& \# \& \# } u_i(b_i, a_{-i}) > u_i(a_i, a_{-i})\}$$

とする。 ρ に関して (a_i, a_{-i}) に対する同等以上の反応対応 (same or better reply correspondence) を

$$R_C(a_i, a_{-i}) = \prod_{i \in N} R_i(a_i, a_{-i})$$

とする。

A の部分集合 A' について、すべての $a \in A$ に対して、 $R_C(a) \cup A'$ なる場合、 A' を同等以上の反応に関して閉じている集合 (Closed Under Same-Or-Better Reply set) とする (以下 CUSOBR 集合)。

A を CUSOBR 集合とする。 R_G の定義から任意の $a \in A$ について $a \in R_G(a)$ である。 したがって $a \in \bigcup_{a \in A} R_G(a)$

である。 また逆に、 任意の $a \in \bigcup_{a \in A} R_G(a)$ に対して、 ある $a' \in A$ が存在して、 $a \in R_G(a')$ である。 CUSOBR 集合の定義から $a \in R_G(a) \subseteq A$ である。 よって、 A が CUSOBR 集合の場合には

$$R_G(A) = \bigcup_{a \in A} \left(\prod_{i \in N} R_i(a_i, a_{-i}) \right)$$

の不動点になっている (つまり $R_G(A) = A$)。

また、 A の部分集合 A' が直積集合で、 かつすべての $a \in A'$ に対して、 $R_G(a) \subseteq A'$ なる場合、 A' を同等以上の反応に関して閉じた行動の組の直積集合 (Product set of action profiles that is Closed Under Same-Or-Better Reply set) とする (以下 PCUSOBR 集合)。 CUSOBR 集合の場合と同様に、 A' が PCUSOBR 集合の場合には

$$R_G(A') = \prod_{i \in N} \left(\bigcup_{a \in A'} R_i(a_i, a_{-i}) \right)$$

の不動点になっている (つまり $R_G(A') = A'$)。

また、 真部分集合となる CUSOBR 集合を持たない CUSOBR 集合 (or PCUSOBR 集合) を最小 CUSOBR 集合 (or 最小 PCUSOBR 集合) とする。

表7 例 (Saran and Serrano [10])

	L	C	R
T	0, 0	0, 0	1, 1
U	1, 0	0, -1	0, 0
M	0, 1	2, 0	0, 0
D	2, 0	0, 1	0, 0

これらの集合を「グラフ」を使って考える。まずゲームの行動の組を頂点とする。いま \mathcal{N} は有限であるから頂点も有限個になる。各頂点と他の頂点を結ぶ辺を考える(こうした頂点と辺からなる構造をグラフという)。いま次のようにグラフの辺に「向き」を付ける(辺に向きの付いたグラフを有向グラフという)。頂点 (a_1, a_2, \dots, a_N) と $(a'_1, a'_2, \dots, a'_N)$ の間には、 $(a_1, a_2, \dots, a_N) \rightarrow (a'_1, a'_2, \dots, a'_N)$ かつ $(a'_1, a'_2, \dots, a'_N) \in R_C(a_1, a_2, \dots, a_N)$ であるとき、この (a_1, a_2, \dots, a_N) から $(a'_1, a'_2, \dots, a'_N)$ へ向かって矢印を付ける。例えば、次の利得表で表される戦略形ゲームを考える (Saran and Serrano [10] Sec.3)。

\mathcal{N} のゲームでは、各 $R_1 (i=1, 2)$ は次のようになる。

- $R_1(T, L) = \{T, U, D\} ; R_2(T, L) = \{L, R\}$
- $R_1(T, C) = \{T, M\} ; R_2(T, C) = \{C, R\}$
- $R_1(T, R) = \{T\} ; R_2(T, R) = \{R\}$
- $R_1(U, L) = \{U, D\} ; R_2(U, L) = \{L\}$
- $R_1(U, C) = \{U, M\} ; R_2(U, C) = \{C\}$
- $R_1(U, R) = \{T, U\} ; R_2(U, R) = \{R\}$
- $R_1(M, L) = \{U, M, D\} ; R_2(M, L) = \{L\}$
- $R_1(M, C) = \{M\} ; R_2(M, C) = \{L, C\}$
- $R_1(M, R) = \{T, M\} ; R_2(M, R) = \{L, R\}$
- $R_1(D, L) = \{D\} ; R_2(D, L) = \{L, C\}$

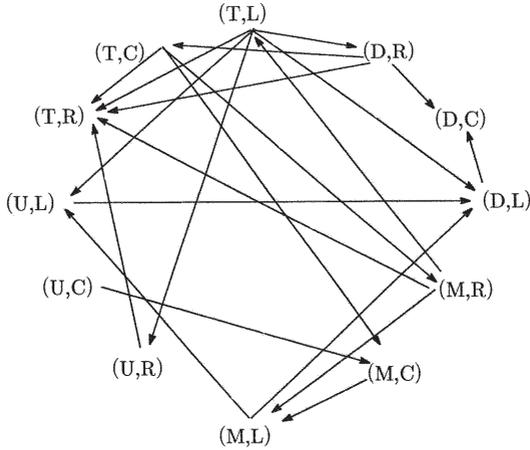


図1 同等反応グラフ

$$R_1(D, C) = \{M, D\} ; R_2(D, C) = \{C\}$$

$$R_1(D, R) = \{T, D\} ; R_2(D, R) = \{C, R\}$$

これから次の R_c 群が得られる。

$$R_c(T, L) = \{(T, L), (T, R), (U, L), (U, R), (D, L), (D, R)\}$$

$$R_c(T, C) = \{(T, C), (T, R), (M, C), (M, R)\}$$

$$R_c(T, R) = \{(T, R)\}$$

$$R_c(U, L) = \{(U, L), (D, L)\}$$

$$R_c(U, C) = \{(U, C), (M, C)\}$$

$$R_c(U, R) = \{(T, R), (U, R)\}$$

$$R_c(M, L) = \{(U, L), (M, L), (D, L)\}$$

$$R_c(M, C) = \{(M, L), (M, C)\}$$

$$R_c(M, R) = \{(T, L), (T, R), (M, L), (M, R)\}$$

$$R_c(D, L) = \{(D, L), (D, C)\}$$

$$R_c(D, C) = \{(M, C), (D, C)\}$$

$$R_c(D, R) = \{(T, C), (T, R), (D, C), (D, R)\}$$

以上から有向グラフをかくと図1のようになる（ここでは自分自身に向いている辺は省いている。以下では同等反応グラフと呼ぶ）。

同等反応グラフで自分以外へ出ていく辺がないものをシンクと呼ぶ。このグラフでは頂点 (T, R) がシンクである。自分以外へは向かわないということは、誰も行動を変更するものがない、つまりシンクは純粹戦略ナッシュ均衡になる。また、頂点 (d, b) (c, c) で $d = T, U, M, D : b = L, C, R$ から (d^{1+}, b^{1+}) へ向けて繋がる辺があるとき (d^1, b^1) は (d^{1+}, b^{1+}) へ隣接するとき (d^1, b^1) 頂点 $(d^1, b^1), (d^2, b^2), \dots, (d^L, b^L)$ が隣接した頂点の列のとき同等反応パスという。図1のグラフから最小CUSOBR集合は (T, R) と $(U, L), (M, L), (M, C), (D, L), (D, C)$ であることがわかり、⁽⁶⁾ また最小PCUSOBR集合は (T, R) だけである。

ある状態の集合 H があり、もし上記の学習プロセスがいったん H のどれかに達したら、 H から抜け出ることがなく、かつそうした性質を満たす真部分集合を含まない場合、 H を再帰クラス (recurrent class) という。

命題 1 (Saran and Serrano [10] Prop.4.1)

任意の状態の集合 $H \subseteq A^m$ に対して、 H のある状態でプレーされる行動の組すべてを $A(H)$ で表すとする。 A がゲームの最小PCUSOBR集合であるとき、 $A(H) \subseteq A$ なるマルコフ・プロセス $M(a)$ の再帰クラス H が存在する。

また、ゲームが同等反応の下で弱非周期的であるとは、いかなる頂点からも少なくとも一つのシンク (ナッシュ均衡) へ向かう同等反応パスがある場合をいう。したがって、もしゲームが同等反応の下で弱非周期的であるなら、すべてのCUSOBR集合は純粹戦略ナッシュ均衡を含む。

これらから次が得られる。

系1 (Saran and Serrano [10] Cor.4.2)

同等反応の下で弱非周期的なゲームを考える。 $A(H) \cup A$ が純粹戦略ナッシュ均衡であることと、 H がマルコフ・プロセス $M(q)$ の再帰クラスになることは同値である。

この系1によって、Hart and Mas-Colell [4] —記憶が有限でない場合—では不明であった点が明らかになる。つまり、期ごとの行動の組 ρ はほとんど確実に有限時間で純粹戦略ナッシュ均衡に達するということである。⁽⁷⁾

四. 結びに代えて

本稿では、学習理論の中でも特に心理学で用いられることの多い「強化学習」、特に「後悔の度合いに応じた行動を選ぶ戦略」を概観した。この分野の理論研究も進んで来ているが、実験で得られた結果をどこまで説明できるかという点に関しては、経済学で使われてきた「信念学習」との比較で、どちらが一般に優れているかといった結論は得られていないようである。

強化・信念学習双方のさらなる研究が望まれるが、いかにシンプルで実験データに適合したモデルを作れるかが問題となる。「後悔しない」行動を取る確率が高くなるという考え方はわかりやすいが、どの程度実用性があるかは今後十分研究・検討されなければならない。本稿後半で、Saran and Serrano [10] による有限の記憶への拡張のケースを見たが、これ以外にも指数関数的に記憶を割り引いて行く方法などが考えられている (Marden, Arslan and Shamma [7])。

参考文献

- [1] Blackwell,D., "A vector valued analog of mini-max theorem", Pacific Journal of Mathematics, 6, 1956, 1-8.
- [2] Fudenberg,D. and Levine,D., "Conditional Universal Consistency", Games and Economic Behavior, 29, 1999, 104-130.
- [3] Hannan,J., "Approximation to Bayes Risk in Repeated Plays", in Contribution to the Theory of Games, 3, 97-139, Princeton Univ. Press, 1957
- [4] Hart,S. and Mas-Colell,A., "A Simple Adaptive Procedure Leading to Correlated Equilibrium", Econometrica 68, 2000, 1127-1150.
- [5] Lehrer, E., "Approachability in infinite dimension space", International Journal of Game Theory 31, 2002, 253-268.
- [6] Lehrer, E., "A wide range no-regret theorem", Games and economic behavior 42, 2003, 101-115.
- [7] Marden,J.R.,Arslan,G. and Shamma,J.S., "Regret based dynamics: Convergence in weakly acyclic games", AAMAS 2007, Proceedings of the 6th international joint conference on Autonomous agents and multiagent systems, 2007.
- [8] Nash,J.F., "Equilibrium Points in N-Person Games", Proceedings of the National Academy of Science of the United States of America 36, 1950, 48-49.
- [9] Nash,J.F., "Non-Cooperative Games", Annals of Mathematics 54, 1951, 286-295.
- [10] Saran, R. and Serrano, Roberto., "Regret matching with finite memory", Working Paper Institute Madriedo de Estudios Aranzados (IMDEA) Ciencias Sociales, 2010.
- [11] 川越敏司 『行動ゲーム理論入門』ニール出版、二〇一〇年。
- [12] van Damme, E., Stability and Perfection of Nash Equilibria, Springer-Verlag,Berlin, 1987.
- [13] Young, H.P., *Nit Strategic Learning and its Limits*, Oxford Press, 2004.

注

- (1) 八十年代くらいまで、ゲーム理論での研究はプレイヤーの合理性を仮定し、ナッシュ均衡をいかに絞り込むかに焦点が置かれていた(ナッシュ均衡の「精緻化」と呼ばれる)。こうした問題に関する一つの到達点が van Damme [12] である。
- (2) 以下では数学的に厳密な議論は行っていない。また、regret matching モデルについて網羅するものでもない。
- (3) ここで集合 R は可測で、 μ を R 上の確率測度とする。
- (4) Blackwell の approachability 定理については Yong [13] 第4章に簡潔な解説がある。
- (5) もし最小 CUSOBR が直積集合であるなら、それは最小 PCUSOBR 集合であり、さらに純粋戦略ナッシュ均衡は一点からなる最小 CUSOBR 集合でもあり、最小 PCUSOBR 集合でもある。また、すべての最小 PCUSOBR 集合はどれかの最小 CUSOBR 集合に含まれる。
- (6) グラフから頂点 (T, R) は「点のみで同等以上の反応に関して「閉じて」おり、 $(U, L), (M, L), (M, C), (D, L), (D, C)$ も「閉じて」いることがわかる。
- (7) この系では「同等反応の下で弱非周期的なゲーム」という条件がつけられているが、Saran and Serrano [10] では、さらに、各プレイヤーが直近 β 期の中からランダムに s 期抜き取り、後悔の度合いを考慮するという場合も分析されている。この場合 β が十分小さければ、 $A(\beta)$ が最小 PCUSOBR 集合であることと、 H がマルコフ・プロセス $M(q)$ の再帰クラスになることとは同値であることが示されている (Saran and Serrano [10] Cor.4.3)。