

模倣ダイナミクスと試行錯誤ダイナミクス

Imitation Dynamics and Trial and Error Dynamics

大 浦 宏 邦
(Hirokuni Oura)

進化ゲーム理論は動学化されたゲーム理論で、遺伝ダイナクスモデルと学習ダイナクスモデルの二階建ての構造を持っている（石原、金井2002）。

学習ダイナクスは、プレーヤーが戦略を何らかの方法で修正することによって生じるダイナクスで、比較的短いタイムスケールの社会現象をモデル化するのに向いている。本論文では学習ダイナクスのうち、模倣ダイナクスと試行錯誤ダイナクスについて、やや詳しく見てみることにしよう。

1 遺伝ダイナクスと学習ダイナクス

まず、遺伝ダイナクスと学習ダイナクスの違いについて確認する。

遺伝ダイナクスモデルの代表的なものはレプリケーターダイナクスモデルだが、これは遺伝や教育で戦略が親から子へ継承される場合に生じるダイナクスである。このモデルではプレーヤーは生涯に渡って戦略を変えることはない。親は子供を残して死んでいくが、この出生死滅過程が

ダイナミクスを駆動するマイクロプロセスとなる。

これに対し、プレーヤーが途中で戦略を変えることを仮定するダイナミクスを考えることもできる。このタイプのダイナミクスを学習ダイナミクスと総称する。現実の社会ではプレーヤーは様々な方法で戦略を変えている。ある戦略で失敗すれば、それを変えるかもしれないし、他人がよさそうな戦略をとっていれば、それをまねするかも知れない。いろいろ考えた結果、現状では最善と判断する戦略をとることもあるだろう。

このように、戦略を変える手続きにはさまざまなものが考えられるので、学習ダイナミクスにも色々なタイプのもものが考えられる。具体的には、試行錯誤ダイナミクス、模倣ダイナミクス、最適反応ダイナミクスなどが考えられている。

以下ではまず、レプリケーターダイナミクスと学習ダイナミクスの特徴をおおまかに比較してから、試行錯誤ダイナミクスと模倣ダイナミクスについて詳しく見ていくことにしよう。

1.1 ゲームの繰り返しと戦略修正

学習型のモデルの最大の特徴は、プレーヤーが戦略を変えることが想定される点である。プレーヤーはゲームを何度か繰り返して経験するうちに、戦略を変えたり変えなかったりする。戦略の変え方には、自分の手や他人の手を参照したり、自分の利得や他人の利得を参照したり、深く考えたり余り考えなかったり様々な方法がありうるが、いずれにせよ何らかの手続きに従ってプレーヤーは戦略を修正していく。この戦略を修正する手続きを戦略修正アルゴリズムという。学習ダイナミクスは、各プレーヤーが何らかの戦略修正アルゴリズムに従って戦略を変化させるときに生じるダイナミクスである。

戦略が修正されるということは、ゲームが繰り返して複数回行われることを意味している。逆にいうとゲームが1回切りしか行われない場合には学習ダイナミクスは生じない。ちなみに従来型のゲーム理論にも「繰り返しゲーム」というカテゴリーがあるが、従来の「繰り返しゲーム」は動学

的なプロセスを扱うものではない。プレイヤーは事前に「繰り返しゲーム」の中でどう振舞うかについての行動プラン（この行動のプランをスーパー戦略という）を持っていると想定され、この行動プラン（スーパー戦略）同士の間でのNash均衡を考えるという静学的なアプローチが取られてきた。「繰り返しゲーム」という用語はこの種の静学的な分析法を採用するときの用語なので、ここでは「繰り返しゲーム」とは思わず、「繰り返しのあるゲーム」「ゲームの繰り返し」などという言い方を使うことにしよう。

1.2 現象のタイムスケール

学習ダイナミクスは、ゲームが複数回繰り返して行われる場合にプレイヤーが戦略を修正することによって生じるダイナミクスである。ダイナミクスの生じるタイムスケールは、ゲームが繰り返される頻度や、戦略の修正が起きる頻度に依存する。ゲームによっては一日に何回かプレーする機会がある場合もあるだろうし、一年に一回とか数年に一回しかプレーする機会がない場合もあるかもしれない。しかし、いずれの場合もプレイヤーが生きているうちに数回以上はゲームを行なう機会と戦略を変更する機会があると考えられるので、長くても1世代以内のタイムスケールで生じるダイナミクスであるといえる。実時間にして数年から数十年程度のタイムスケールであろう。

これに対し、レプリケーターダイナミクスの場合は出生死滅過程によってダイナミクスが駆動する。この場合に想定されるタイムスケールは数十世代から数百世代で、実時間にして数百年から数千年以上のタイムスケールになる。これらより、学習ダイナミクスによる分析は、短期的な現象を分析する場合に有用であり、レプリケーターダイナミクスによる分析は中長期的な現象を分析する場合に有用であるといえる。

生物学的な現象の場合は、数千年から数万年というタイムスケールで起こることも多いが、社会科学の分野では数年から数十年程度の現象がほとんどであろう。したがって、多くの社会科学的な現象の分析には学習ダイナミクスが用いられることになる。しかし、社会科学的な現象でも社会の

起源や協力の進化といった問題にかかわる現象は、長いタイムスケールで起きる現象なのでレプリケーターダイナミクスによる分析が有用であり必要となる。

以下では、具体的な学習ダイナミクスとして、模倣ダイナミクスと試行錯誤ダイナミクスについて見てみよう。

2 模倣ダイナミクス

模倣ダイナミクスは、プレーヤーが他者の戦略をまねることによって生じるダイナミクスである。ここでは2人2戦略の対称ゲーム（自分と相手の立場を入れ替えても戦略や利得が変わらないゲーム）の場合について模倣ダイナミクスのモデルを紹介する。

2.1 模倣と同調

模倣とは、簡単にいうと他のプレーヤーの戦略をまねすることである。他人の服装や仕草をまねて取り入れることは良くあるし、ものの考え方や生き方をまねるといふこともあるだろう。

模倣と良く似た現象に同調がある。他者をまねるといふ点で同調と模倣は共通している。相違点は、同調が多くの人が行っている行動を「みんなしてるから」という理由でまねるのに対し、模倣では他者の行動などを「いいな」と思ってまねる点にある。同調では周りと同じかどうかの方が大事であって、その行動そのものが良いかどうかはあまり問題にならないのに対し、模倣ではその行動が良いか悪いかの方が問題で、良ければまねするし悪ければまねしない（その結果相手と違う行動となってもかまわない）点が異なっているといえる。

もちろん実際の行動はこれらの要素が混在しているが、理念型としては模倣と同調を上のように区別することができる。同調によるダイナミクスは「多い方の戦略が増える」という単純なものとなりやすい。以下では、他人の行動などを「いいな」と思ってまねをする場合のダイナミクスを考

えることにする。

2.2 仮定すべき事柄

模倣によってプレーヤーの戦略は変化し、それにつれて集団中の戦略シェアも変化していく。このプロセスをモデル化する上で、考慮すべき点を確認しておこう。

まず、一口に「他人の行動」をまねするといっても、具体的に誰の行動をまねすると想定すればよいのであろうか。模倣対象についての仮定が必要である。次に「いいな」と思う場合の「良い悪い」を判断する基準はなんなのであろうか。評価基準についての仮定が必要である。さらに、他人の何をまねするのか（模倣事項）、どういうときにまねするのか（模倣タイミング）についての仮定も必要であろう。

模倣対象については、集団中からランダムに選ぶ（たまたま道で会った人を対象とするなど）と想定することもできるし、羽振りのよさそうな人を対象にすると想定することもできる。

良い悪いの評価基準としては、効用関数を想定することが一般的である。効用の高い行動を良い、低い行動を悪いと評価するのである。ただしこの場合、模倣しようとする人の効用関数が、模倣される人の効用関数が、どちらが妥当かということが問題になる。誰かが気にいって着ている服でも、他の人にはちっとも良くないということがありうる。この場合、良くないと思っている人はその服装をまねしないであろう。したがって、評価の基準は「模倣しようとする人の効用関数」とすることが妥当である。

他人の何をまねするのか（模倣事項）については経験的な研究が必要であるが、次の事柄が指摘できるであろう。まず、他人から見えにくい事柄より見えやすい事柄の方が模倣の対象となりやすいと考えられる。例えば下着や寝巻きよりも上着やコートの方が模倣されやすいし、流行の対象にもなりやすいであろう。また、複雑な事柄よりも単純な事柄の方が模倣されやすいであろう。複雑な事柄が模倣されるときには、単純化されて模倣される可能性も考えられる。その意味で、混合戦略（複数の戦略を確率的

に組み合わせる)やスーパー戦略(相手の出方に応じた行動プラン)は模倣されにくく、模倣される場合もそのままではなく純粹戦略として模倣されることが多いと予想される。ここでは、単純な純粹戦略が模倣の対象となる場合を考える。

模倣のタイミングについては、すべてのプレーヤーが同じ頻度で他人の行動を参照しようとする想定することもできるし、特に現状に不満をもっている人が高い頻度で他人の行動をまねしようとすることもあるであろう。この点については、両方のタイプのモデルを立てることができる。

2.3 街角模倣モデル

ここではまず、もっとも単純なモデルとして次の事柄を仮定するモデルを考えよう。

[仮定]

- ・どのプレーヤーも同じ頻度で戦略の見直しをする。
- ・効用の比較対象は集団中からランダムに選ぶ。
- ・各プレーヤーは微小時間 dt の間に $r dt$ の確率で戦略の見直しをする。
- ・プレーヤー集団全体の人数は N 人で一定。

効用の比較対象を集団中からランダムに選ぶということは、街角で偶然出会った人の戦略をまねるかまねないか意思決定するようなモデルなので、このタイプのモデルを街角模倣モデルと呼ぶ。

各プレーヤーは dt の間に $r dt$ の確率で戦略の見直しをするので dt の間に $r N dt$ 人のプレーヤーが戦略の見直しを行う(N は十分に大きいものとする)。

見直しを行うプレーヤーは集団中からランダムに一人、参照相手となるプレーヤーを選び出す。参照相手が自分よりも自分の効用関数でみて良い場合に、自分の戦略を参照相手の戦略に変更するものとする。ただし、以下では効用関数はどのプレーヤーについても共通であると仮定する。この

ようなマイクロプロセスを仮定するときに、集団全体ではどのようなマクロダイナミクスが生じるのであろうか。いくつかの具体的なゲームについて見てみよう。

[例] 調整ゲーム

次の調整ゲームの場合について考えよう。

自分 \ 相手		A	B
A		2	0
B		0	1

戦略Aのプレーヤーの割合を x 、戦略Bのプレーヤーの割合を $1 - x$ とする。プレーヤーはランダムに出会ってゲームを行うとすると、最近のプレーで効用2を獲得したプレーヤーの割合は x^2 、効用1を獲得したプレーヤーの割合は $2x(1 - x)$ 、効用0を獲得したプレーヤーの割合は $(1 - x)^2$ となる。

微小時間 dt の間に $r N dt$ 人のプレーヤーが戦略の見直しを行うが、そのうち

- $x^2 r N dt$ 人が 最近の効用が 2
- $2x(1 - x) r N dt$ 人が 最近の効用が 1
- $(1 - x)^2 r N dt$ 人が 最近の効用が 0

となる。ここで各プレーヤーは、参照相手の最近のゲーム結果が分かり、それと自分の最近のゲーム結果を比較して戦略を変更するかしないかを定めることとしよう。

このとき、戦略Aのプレーヤーが戦略Bに変更するのは、戦略Aで効用0のプレーヤーが、戦略Bで効用1のプレーヤーを参照相手に選んだ場合である。このときAからBへの「流出」が発生する。

見直しプレーヤー $r N dt$ 人のうち

- ・戦略Aで効用0のものは $x(1 - x) r N dt$ 人

- ・これらのプレーヤーが戦略Bで効用1のプレーヤーを参照相手に選ぶ確率は $(1 - x)$

である。したがって

$$x(1 - x)r Ndt \cdot (1 - x)人$$

がdtの間にAからBに流出する人数となる。

一方、戦略Bから戦略Aに変更するのは、戦略がBで効用が0や1だったプレーヤーが戦略Aで効用が2だったプレーヤーを参照する場合である。この場合、戦略Aへの流入が発生する。

戦略Bのプレーヤーは全員効用が0か1なので、 $r Ndt$ 人の見直しプレーヤーのうち、

$$(1 - x)r Ndt \cdot x^2人$$

のプレーヤーが、dtの間に戦略をBからAに変更することになる。これがAへの「流入」である。

ここで戦略Aの人数は、Aへの流入人数からAからの流出人数を引いた人数だけ増減するはずである。したがって、

$$Aの人数の変化 = Aへの流入 - Aからの流出$$

となる。上の考察から

$$Aへの流入 = (1 - x)r Ndt \cdot x^2$$

$$Aへの流出 = x(1 - x)r Ndt \cdot (1 - x)$$

なので

$$Aの人数の変化 = (1 - x)r Ndt \cdot x^2 - x(1 - x)r Ndt \cdot (1 - x)$$

である。

また、戦略Aのシェアxの変化をdxとすると

$$dx = Aの人数の変化 / N$$

である。したがって

$$Aの人数の変化 = Ndx$$

となる。これを上の式に代入すると

$$Ndx = (1 - x)r Ndt \cdot x^2 - x(1 - x)r Ndt \cdot (1 - x)$$

両辺をNで割って

$$dx = (1 - x)r dt \cdot x^2 - x(1 - x)r dt \cdot (1 - x)^2$$

両辺をdtで割って

$$\begin{aligned} dx/dt &= (1 - x)r \cdot x^2 - x(1 - x)r \cdot (1 - x)^2 \\ &= r x(1 - x) \{ x - (1 - x)^2 \} \\ &= r x(x - 1) \{ x^2 - 3x + 1 \} \end{aligned}$$

という微分方程式が得られる。これが街角模倣モデルのダイナミクス方程式である。

このダイナミクスを分析してみよう。2次方程式の解の公式を用いると

$$x^2 - 3x + 1 = (x - \frac{3 - \sqrt{5}}{2})(x - \frac{3 + \sqrt{5}}{2})$$

$$\text{ただし } \frac{3 - \sqrt{5}}{2} \approx 0.38$$

$$\frac{3 + \sqrt{5}}{2} \approx 2.62$$

と因数分解できる。これよりダイナミクスのベクトル図は図1のようになる。

したがって、

$$x = 0, x = 1 \text{ が 漸近安定点}$$

$$x = \frac{3 - \sqrt{5}}{2} \approx 0.38 \text{ が 不安定定常点}$$

であることがわかる。

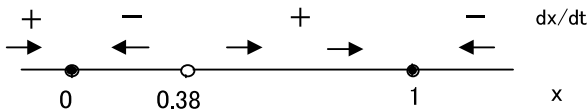


図1 街角モデルのダイナミクス

ちなみに、利得表の値を効用ではなく適応度の増減とみなしてレプリケーターダイナミクスを求めると

$$dx/dt = x(1 - x) \{ 3x - 1 \}$$

となり

$$x = 0, 1 \text{ が 漸近安定点}$$

$$x = 1/3 \text{ が 不安定定常点}$$

となる。

この場合、レプリケーターダイナミクスと模倣ダイナミクスはおおむね同じようなダイナミクスとなる。不安定定常点（ダイナミクスの分水嶺にあたる）の位置は若干異なるが、漸近安定状態については両者は一致している。

2.4 不満・羽振りモデル

街角模倣モデルでは、見直しの頻度はすべてのプレーヤーで同じで、模倣対象者は集団からランダムに選ばれた。しかし、戦略の見直しは効用の低い人が頻繁に行い、効用の高い人はあまり行わないかもしれない。プレーヤーが自らの効用の低さに不満をもち、頻繁に他者の戦略を参照しようとすることを仮定するモデルを不満モデルと呼ぶことにする。

このとき、参照者も集団からランダムに等確率で選ばれるのではなく、「羽振り」の良い人がそうでない人よりも高い確率で選ばれるかもしれない。この場合、プレーヤーの効用が観察可能な経済的・社会的な「羽振り」と関連していることを仮定する必要があるが、なんらかのメカニズムで羽振りの良いプレーヤーが高い確率で参照されることを仮定するモデルを羽振りモデルと呼ぶことにする。

以下では、不満による模倣と羽振りのよさによる参照の両方を仮定する不満・羽振りモデルを考えることにしよう。

〔仮定の追加〕

以下のモデルでは、街角モデルに次の仮定を追加する。

- ・最近の効用が u であったプレーヤーは、微小時間 dt の間に $r(u)dt$ の確率で戦略の見直しを行う。
- ・最近の効用が u_1 であったプレーヤーは、効用が u_2 であったプレーヤーより $q(u_1)/q(u_2)$ 倍参照されやすい。
- ・ $r(u)$ は u についての減少関数。 $q(u)$ は u についての増加関数。

まず、見直しの仮定である。最近のプレーで効用が u であったプレーヤーは、微小時間 dt の間に $r(u)dt$ の確率で戦略の見直しを行うものとする。 $r(u)$ は u についての減少関数だが、負にはならないものとしよう。このように仮定すれば効用が低いプレーヤーほど頻繁に見直しを行うことになる。

次に参照についての仮定である。最近のプレーで効用が u_1 であったプレーヤーは、効用が u_2 であったプレーヤーより $q(u_1)/q(u_2)$ 倍参照されやすいと仮定する。ここで $q(u)$ は u についての増加関数で、常に正の値をとるものとする。このように仮定すると効用の高いプレーヤーほど参照相手として選ばれやすいことになる。

たとえば、仮に $q(1) = 1$ 、 $q(2) = 2$ 、だとするならば、効用が 2 のプレーヤーは効用が 1 のプレーヤーよりも 2 倍の確率で他のプレーヤーに参照されることになる。もし、すべてのプレーヤーの効用が 1 か 2 で、効用 1 のプレーヤーの割合が x_1 、効用 2 のプレーヤーの割合が x_2 ならば、誰か一人参照相手のプレーヤーを選んだときにそのプレーヤーの効用が 1 である確率は

$$x_1 / (x_1 + 2x_2)$$

となる。また参照プレーヤーの効用が 2 である確率は

$$2x_2 / (x_1 + 2x_2)$$

となる。

一般には、効用 u_i のプレーヤーの割合が x_i で、参照されやすさが $q(u_i)$ に比例するとき、一人選んだ参照プレーヤーの効用が u_j である確率は

$$q(u_j)x_j / \sum_i q(u_i)x_i$$

となる。

【例】 調整ゲーム

再び次の調整ゲームについて考えよう。不満・羽振りモデルの場合、どのようなダイナミクスが生じるのであろうか。

自分\相手	A	B
A	2	0
B	0	1

まず、戦略Aのシェアを x 、戦略Bのシェアを $1 - x$ とし、ランダムマッチングでゲームが行われているものとする、最近のプレーで

効用2のプレーヤーの割合 x^2

効用1のプレーヤーの割合 $(1 - x)^2$

効用0のプレーヤーの割合 $2x(1 - x)$

となる。

微小時間 dt の間に効用2のプレーヤーが見直しをする確率は $r(2)dt$ なので、集団全体 N 人のうちでは $x^2N r(2)dt$ 人の効用2のプレーヤーが戦略の見直しを行うことになる(N は十分に大きいものとする)。効用1や効用0のプレーヤーについても同様である。

ここで、 dt の間に戦略をAからBに変更する人数を求めることにしよう。AからBに変更するのは、戦略がAで効用が0であったプレーヤーが、戦略がBで効用が1であったプレーヤーを参照する場合である。

dt の間に戦略Aで効用0のプレーヤーのうち

$$x(1 - x)N r(0)dt \text{人}$$

が戦略の見直しを行う。これらのプレーヤーは集団から参照相手のプレーヤーを一人選び出すが、そのプレーヤーが戦略Bで効用1である確率は

$$q(1)(1 - x) / (2q(0)x(1 - x) + q(1)(1 - x) + q(2)x^2)$$

である。以下、分母を Q と書くことにすれば、確率は

$$q(1)(1 - x) / Q$$

となる。これより、戦略A効用0で戦略B効用1のプレーヤーを参照して戦略をBに変更する人数は

$$x(1 - x)N r(0)dt \cdot q(1)(1 - x) / Q \text{人}$$

となる。これが、 dt の間にAからBに流出する人数である。

次に、BからAに流入する人数を考える。戦略Bから戦略Aに変更する

ケースは、戦略Bで効用0のプレーヤーが戦略Aで効用2のプレーヤーを参照する場合と、戦略Bで効用1のプレーヤーが戦略Aで効用2のプレーヤーを参照する場合の2通りある。

流出の場合と同様に、それぞれの人数を求めると、前者のケースでの流入は

$$x(1-x)Nr(0)dt \cdot q(2)x^2/Q人$$

であり、後者のケースでの流入は

$$(1-x)Nr(1)dt \cdot q(2)x^2/Q人$$

となる。これらの合計が微小時間dtの間に戦略Aに流入する人数となる。

流入と流出の人数がわかれば、街角モデルと同様にAの割合xについての微分方程式を立てることが出来る。

$$(Aの増減) = (Aへの流入) - (Aからの流出)$$

であるから

$$\begin{aligned} Ndx = & x(1-x)Nr(0)dt \cdot q(2)x^2/Q \\ & + (1-x)Nr(1)dt \cdot q(2)x^2/Q \\ & - x(1-x)Nr(0)dt \cdot q(1)(1-x) \end{aligned}$$

両辺をNdtで割ると

$$\begin{aligned} dx/dt = & x(1-x)r(0) \cdot q(2)x^2 + (1-x)Nr(1) \cdot q(2)x^2 \\ & - x(1-x)r(0) \cdot q(1)(1-x) \\ = & x(1-x)r(0) \cdot q(2)x^2 + (1-x)r(1) \cdot q(2)x \\ & - r(0) \cdot q(1)(1-x) \\ = & x(1-x)r(0)q(2) - r(0)q(1) - r(1)q(2)x^2 \\ & + (2r(0)q(1) + r(1)q(2))x - r(0)q(1) \end{aligned}$$

となる。これが不満・羽振りモデルにおける模倣ダイナミクス方程式である。

これは、どのようなダイナミクスになるのであろうか。このままでは、分かりにくいので仮に

$$r(0)=0.5 \quad r(1)=0.3 \quad r(2)=0.1$$

$$q(0)=1 \quad q(1)=3 \quad q(2)=5$$

として計算してみると

$$dx/dt = x(x - 1)(x^2 + 9x - 3)/2Q$$

となる。

この場合のベクトル図は図2のようになり、

$x = 0, 1$ が 漸近安定点

$x = 0.32$ が 不安定定常点

ただし $x = (-9 + \sqrt{93})/2$ (約0.32)

となる。街角モデルの場合と不安定定常点の位置がやや異なる(街角モデルでは約0.38、不満・羽振りモデルでは約0.32)が、定性的にはほぼ同じダイナミクスとなることが分かる。

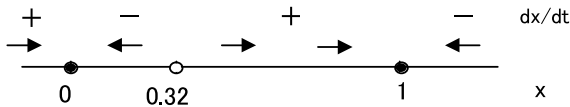


図2 不満・羽振りモデルのダイナミクス

2.5 2 × 2 ゲームの模倣ダイナミクス

街角モデル、不満・羽振りモデルのいずれでも調整ゲームの漸近安定状態は $x = 0$ と $x = 1$ であった。調整ゲームは (A, A) と (B, B) が strict Nash 均衡となるゲームだが、この結果から模倣ダイナミクスでは strict Nash 均衡に対応する状態が漸近安定になることが予想される。ここでは、一般の 2 人 2 戦略対称ゲーム (2 × 2 対称ゲーム) の模倣ダイナミクスをたてることでこの予想を確認しておこう。

自分 \ 相手	A	B
A	a	b
B	c	d

という 2×2 ゲームを考える。モデルの仮定は基本的に不満・羽振りモデルと同じだが、さらに次の仮定を付け加える。

[追加の仮定]

- ・ 効用 u_1 のプレーヤーが効用 u_2 のプレーヤーを参照したときに、
 戦略を効用 u_2 の戦略に変更する確率を $(u_2 - u_1)$ とする。
 (ただし、 α は 0 から 1 までの値を取る $u_2 - u_1$ に関する増加関数)

先ほどまでのモデルでは、参照相手の効用が高ければ相手の戦略を採用し、低ければ採用しないことを仮定していたので、戦略変更確率 $(u_2 - u_1)^\alpha$ が

$$u_2 > u_1 \quad \text{のとき} \quad (u_2 - u_1)^\alpha = 1$$

$$u_2 < u_1 \quad \text{のとき} \quad (u_2 - u_1)^\alpha = 0$$

となる場合を考えていたことになる。

これをもう少し一般化して、相手との差が小さいときは戦略をまねするかどうか微妙だが、相手の効用の方がずっと良いときには高い確率でまねをすることを想定したのが、ここで追加をした仮定である。参照相手の効用を見極める上である程度の誤差がある場合や、相手のまねをすることにある程度のためらいがある場合には、この仮定のように考えると考えられる。

以下、このような場合の模倣ダイナミクスを考えてみよう。

戦略 A のシェアを x 、戦略 B のシェアを y とし、ランダムマッチングでゲームが行われているものとする、最近のプレーで

$$\text{効用 a のプレーヤーの割合} \quad x^2$$

$$\text{効用 b のプレーヤーの割合} \quad x y$$

$$\text{効用 c のプレーヤーの割合} \quad x y$$

$$\text{効用 d のプレーヤーの割合} \quad y^2$$

となる。

微小時間 dt の間の見直しで戦略 A から戦略 B への変更が生じるのは、

- 1) 効用 a のプレーヤーが効用 c のプレーヤーを参照した場合

- 2) 効用 a のプレーヤーが効用 d のプレーヤーを参照した場合
- 3) 効用 b のプレーヤーが効用 c のプレーヤーを参照した場合
- 4) 効用 b のプレーヤーが効用 d のプレーヤーを参照した場合

の 4 通りである。

1) のケースで変更を行なう人数は全体の人数 N が十分に多いとして

$$1) \quad r(a)x^2 \cdot q(c)xy / Q \cdot (c - a)N \text{人}$$

である。(ただし $Q = q(a)x^2 + q(b)xy + q(c)xy + q(d)y^2$)

同様に、2) 3) 4) のケースで戦略を A から B に変更する人数はそれぞれ

$$2) \quad r(a)x^2 \cdot q(d)y^2 / Q \cdot (d - a)N \text{人}$$

$$3) \quad r(b)xy \cdot q(c)xy / Q \cdot (c - b)N \text{人}$$

$$4) \quad r(b)xy \cdot q(d)y^2 / Q \cdot (d - b)N \text{人}$$

である。これらの合計が、 dt の間の A からの流出人数である。

次に A への流入人数を考える。 dt の間に A への流入が生じるのは

$$5) \quad \text{効用 c のプレーヤーが効用 a のプレーヤーを参照した場合}$$

$$6) \quad \text{効用 c のプレーヤーが効用 b のプレーヤーを参照した場合}$$

$$7) \quad \text{効用 d のプレーヤーが効用 a のプレーヤーを参照した場合}$$

$$8) \quad \text{効用 d のプレーヤーが効用 b のプレーヤーを参照した場合}$$

の 4 通りである。

流出の場合と同様に計算すると、それぞれの場合の流入人数は

$$5) \quad r(c)xy \cdot q(a)x^2 / Q \cdot (a - c)N \text{人}$$

$$6) \quad r(c)xy \cdot q(b)xy / Q \cdot (b - c)N \text{人}$$

$$7) \quad r(d)y^2 \cdot q(a)x^2 / Q \cdot (a - d)N \text{人}$$

$$8) \quad r(d)y^2 \cdot q(b)xy / Q \cdot (b - d)N \text{人}$$

となる。これらの合計が dt の間の A への流入人数である。

以上を

$$(A \text{ の増減}) = (A \text{ への流入}) - (A \text{ からの流出})$$

に代入して、両辺を Ndt で割って整理すれば模倣ダイナミクス方程式が得られる。煩雑になるので途中の計算は省略して結果だけ示せば

$$dx/dt = x y f(x,y)/Q \quad (\text{式 1})$$

ただし、

$$\begin{aligned} f(x,y) = & x^2(r(c)q(a) - (a-c) - r(a)q(c) - (c-a)) \\ & + x y(r(c)q(b) - (b-c) - r(b)q(c) - (c-b)) \\ & + x y(r(d)q(a) - (a-d) - r(a)q(d) - (d-a)) \\ & + y^2(r(d)q(b) - (b-d) - r(b)q(d) - (d-b)) \end{aligned}$$

となる。これがこの場合の模倣ダイナミクス方程式である。

これは

$$u_2 > u_1 \quad \text{のとき} \quad (u_2 - u_1) = 1$$

$$u_2 < u_1 \quad \text{のとき} \quad (u_2 - u_1) = 0$$

とすると、先の不満・羽振りダイナミクスに一致するし、 $r(u)$ 、 $q(u)$ を一定とすると始めの街角モデルに一致する。その意味では、このダイナミクスを一般模倣ダイナミクスと呼ぶことができる。

2.6 頂点の安定性

$f(x,y)$ が煩雑であるが、式1の形自体は簡単である。この式からまず $x = 0$ と $y = 0$ の時は模倣ダイナミクスは定常であることが分かる。つまり、状態空間の両端（頂点）は定常である。

では、頂点の安定性はどうであろうか。これについては次の命題が成立する。

【命題1】 2×2 対称ゲームの一般模倣ダイナミクスでは、戦略プロファイル (I, I) が strict Nash 均衡ならば、頂点 I は漸近安定である。

【証明】

戦略プロファイル (A, A) が strict Nash 均衡であるとする。このとき $a > c$ である。

式1の $f(x,y)$ に $x = 1$ 、 $y = 0$ を代入すると

$$f(1,0) = r(c)q(a)(a-c) - r(a)q(c)(c-a)$$

である。

r は負ではない減少関数、 q は負ではない増加関数なので、 $a > c$ のとき

$$r(c) > r(a) \quad 0$$

$$q(a) > q(c) \quad 0$$

である。

また、 $a > c$ より

$$2a > 2c$$

$$a - c > c - a$$

であり、さらに q は負ではない増加関数なので

$$(a - c) > (c - a) \quad 0$$

がいえる。

これらより、

$$r(c)q(a)(a-c) > r(a)q(c)(c-a)$$

なので

$$f(1,0) > 0$$

である。

$f(x,y)$ は連続な関数なので、 $x = 1$ 、 $y = 0$ に十分近い x 、 y の範囲では

$$f(x,y) > 0$$

が成立する。

これより、 $x = 1$ から x が少し小さくなくてもそのずれが十分に小さければ

$$dx/dt > 0$$

となり、 x は 1 に収束する。したがって頂点 A に対応する $x = 1$ は漸近安定である。

戦略プロファイル (B, B) が strict Nash 均衡の場合も同様に証明できる。

【証明終わり】

この命題は、strict Nash均衡に対応する頂点が漸近安定であるという、先ほどの予測を確認するものである。これは、一般模倣ダイナミクスに関する命題なので、街角モデル、不満・羽振りモデルを含む広いクラスの模倣ダイナミクスについて成立する。

命題 1 に関連して、次の命題が成立する。

【命題 2】 2×2 対称ゲームの一般模倣ダイナミクスでは、戦略プロファイル (I, I) が Nash 均衡でなければ、頂点 I は不安定である。

【証明】

戦略プロファイル (A, A) が Nash 均衡でないとする。このとき $a < c$ である。

命題 1 と同様に

$$f(1,0) = r(c)q(a)(a-c) - r(a)q(c)(c-a)$$

であるが、 $a < c$ のときは

$$r(a) > r(c) \quad 0$$

$$q(c) > q(a) \quad 0$$

$$(c-a) > (a-c) \quad 0$$

である。したがって

$$f(1,0) < 0$$

$f(x,y)$ は連続な関数なので、 $x = 1$ 、 $y = 0$ に十分近い x 、 y の範囲では

$$f(x,y) < 0$$

が成立する。

これより $x = 1$ は不安定である。

戦略プロファイル (B, B) が Nash 均衡でない場合も同様に証明できる。

【証明終わり】

この頂点の安定、不安定に関する命題はレプリケータダイナミクスの場合と共通している。この点で模倣ダイナミクスはレプリケータダイナミクスと良く似たダイナミクスである。

例えば、囚人のジレンマゲームでは、(C、C) はNash均衡ではなく (D、D) がstrict Nash均衡となるため、一般模倣ダイナミクスでは頂点Dが漸近安定となる。

一方、チキンゲームでは、(A、A) も (B、B) もNash均衡ではないので、頂点Aも頂点Bもどちらも不安定である。この場合は内点に漸近安定点を持つことになる。この漸近安定点の位置は r や q や の値に依存するので、レプリケータダイナミクスの内点漸近安定点に一致する保障はない。この点で、模倣ダイナミクスはレプリケータダイナミクスと異なっているといえる。

3 試行錯誤ダイナミクス

この節では、プレーヤーが試行錯誤で戦略を変更する場合を考えてみよう。試行錯誤による学習は、他の方法による学習に比べて必要とする情報量が少ないので、情報が少ない状況で進行する学習をモデル化するのに適している。ここでは、試行錯誤学習を仮定するモデルの代表としてロスとエレブのモデル (Roth & Erev 1995、Erev & Roth 1998) を紹介する。

3.1 試行錯誤による学習

食事に行くときに「あの店はおいしかったのでまた行こう」とか「あそこはいまいちだったからやめておこう」といった形で意思決定をすることは、日常よく経験する事柄である。このように、複数の選択肢があるときに、とりあえず何か試してみても結果の良かった選択肢を高い確率で採用し、結果の悪かった選択肢を避けるタイプの学習を試行錯誤という。模倣と並んで広く用いられている学習の方略である。心理学の分野では強化学習と呼ばれることもある。

このアルゴリズムで学習をする場合、プレーヤーに必要とされる事柄は、行動の選択肢を幾つかもっていることと、行動の結果を評価する評価基準を持っていることの二つである。ある選択肢を採ったときにどうなるかを事前に知っておく必要はないし、他のプレーヤーの戦略を知る必要もない。また、他のプレーヤーの行動結果について知る必要もない。もちろん、これらの情報があればさらに効率よく学習を行なうことができるが、試行錯誤学習はこれらの情報が手に入らないときにも使用することのできる、汎用性の高い学習アルゴリズムということが出来る。

やや余談であるが、学習という現象は植物よりも動物において一般的にみられる。動物は移動して危険を避けつつえさを採るという生活様式を採用しているが、危険のありかもえさのありかも常に同じ場所とは限らず、時々刻々変化するのが普通である。この状況で生活するには、危険やえさのありかに応じて行動を変えることができたほうが有利である。学習という現象は、このような必要性を背景として進化してきたと考えられる。

試行錯誤学習は学習のなかでももっとも基本的な学習で、良い結果（えさがあった、危険を避けられた、など）をもたらす行動を強化し、悪い結果（えさがなかった、危険な目にあったなど）をもたらす行動を抑制することが原型となっている。もう少し具体的には、神経系の中に結果の良し悪しを判定する部門があって、「良い」結果の場合にはその行動をもたらした神経細胞の結合（シナプス結合）を強化する信号を発し、「悪い」結果の場合には悪い結果をもたらしたシナプス結合を弱める信号を発して、神経結合が作り変えられていく。このようにして、「良い」結果をもたらす行動が発生しやすいように神経結合が作り変えられていくプロセスが試行錯誤学習で、人間に限らず動物が一般に行なっているタイプの学習である。

人間の場合は、「他人の経験から学ぶ」ことも多いが、「自分の経験から学ぶ」ことも依然として大きな位置を占めている。ここでは試行錯誤学習がどのようなダイナミクスをもたらすかを考えてみよう。

3.2 強化と忘却

ここでは、試行錯誤学習をモデル化する方法について考えてみよう。

まず試行錯誤の過程で変化していくものは何なのであろうか。この場合、表面的にはプレイヤーの行動が変化しているように見えるが、行動自体は変化していなくても、内心では「この店おいしくなくなってきたな」とか「あの携帯いいかも」といった評価の変化が生じている可能性がある。この内心の変化が大きくなると、実際の行動も変化すると考えられる。ロスとエレブはこの「内心の変化」を差分方程式の形で表現するモデルを考案した。

このモデルでは、プレイヤーはある戦略を採ろうとする「傾向」(propensity)を持つものとする。この「傾向」が大きい戦略は高い確率で採用されるが、「傾向」の低い戦略は低い確率でしか採用されない。上の例でいう「内心の評価」がこの「傾向」に相当する。

プレイヤーがある戦略を実際に採ってみて、結果が良かったときには、その戦略に対する評価が上がり、その戦略を採ろうとする「傾向」が大きくなる。結果が悪かったときには、評価が下がり、その戦略を採ろうとする「傾向」が小さくなる。この効果は学習心理学の用語で強化(reinforcement)と呼ばれる。一般に高い利得は大きな強化をもたらし、低い利得や負の利得は小さな強化やマイナスの強化をもたらすことが知られている。

ただし、一度強化を受けても時間が経つと強化の効果は次第に失われていく。一度おいしいものを食べてもその印象は次第に薄れていくし、かつて嫌な体験をしてもその記憶は次第に薄れていくであろう。これが学習心理学でいう忘却の過程である。忘却についても様々な実験が行われているが、初期の忘却の速度は速く、時間が経つにつれて速度が遅くなるのが一般的な傾向のようである。

この強化と忘却という過程を取り入れて、ある戦略を採ろうとする「傾向」の変化を次のようにモデル化しよう。 $p_{ij}(t)$ をプレイヤー i が時刻 t に戦略 j を採ろうとする傾向の大きさ、 λ_j を忘却の速さを表すパラメー

ター（ただし $0 < \lambda < 1$ ）、 $R_{ij}(t)$ を時刻 t に j に与えられる強化の大きさとして、時刻 $t + 1$ における傾向の大きさ $p_{ij}(t+1)$ が

$$p_{ij}(t+1) = (1 - \lambda)p_{ij}(t) + \lambda R_{ij}(t) \quad (\text{式 2})$$

となると考える。これがロス・エレブモデルの基本方程式である。

式 2 で忘却が無い場合を考えると、傾向 $p_{ij}(t)$ は強化 $R_{ij}(t)$ の分だけ増減する。また、強化が無い場合を考えると傾向 $p_{ij}(t)$ は每期 $(1 - \lambda)$ の割合で減少していく。このとき初期の減少は比較的すみやかで、時間が経つにつれて減少のしかたは次第に緩やかになっていく。このような性質は、学習心理学で得られた知見と少なくとも定性的に一致しているといえる。

3.3 傾向性と行動

次に、心理的傾向と実際の行動との関係を考えておこう。A を採りたくて B を採りたくない人は大体 A を採るであろうし、逆の人は大体 B を採るであろう。しかし、場合によっては A も採りたいし B も採りたいというケースや、A も B も採りたくないというケースも考えられる。これらの場合、プレーヤーは迷ったあげく、A または B をとる、あるいはどちらも採らないという行動をとるであろう。

どちらもとらないという選択肢がないものとする、プレーヤーは A または B を確率的にとると考えるのが、妥当な推論となる。A または B をとりたいと思う傾向が同程度ならば、A や B をとる確率は等しいと考えられるし、どちらかを取りたい傾向がやや大きい場合は、そちらをとる確率がやや大きくなるであろう。この状況を近似的にモデル化するために、「プレーヤーは戦略 i をとりたいと思う傾向に比例した確率で戦略 i をとる」という仮定がおかれることが多い。

この仮定のもとで、傾向性と行動の関係がどうなるかを考えておこう。まず 2 戦略の場合を考える。時刻 t でプレーヤー i が戦略 A をとりたいと思う傾向を $p_{ia}(t)$ 、プレーヤー i が戦略 B をとりたいと思う傾向を $p_{ib}(t)$ とするとき、プレーヤー i が A を取る確率 $x_{ia}(t)$ と B を取る確率 $x_{ib}(t)$ はそれぞれ

$$x_{ia}(t) = p_{ia}(t) / (p_{ia}(t) + p_{ib}(t))$$

$$x_{ib}(t) = p_{ib}(t) / (p_{ia}(t) + p_{ib}(t)) \quad (\text{式 3})$$

となる。 $x_{ij}(t)$ は $p_{ij}(t)$ に比例しているし、 $x_{ia}(t) + x_{ib}(t) = 1$ となり確率の要件を満たしている。

同様に、戦略が1からnまでn個ある場合は

$$x_{ij}(t) = p_{ij}(t) / p_{ik}(t)$$

となる。やはり、 $x_{ij}(t)$ は $p_{ij}(t)$ に比例しているし、 $x_{ij}(t) = 1$ となり確率の要件を満たしている。

ただしこの場合、 $p_{ij}(t) < 0$ となる場合があると確率が負になってしまつて都合が悪い。そこで式2で、強化 $R_{ij}(t)$ は正または0だと考えるか、負の強化がある場合は $p_{ij}(t) = 0$ となると、それ以上 $p_{ij}(t)$ は減らないことを仮定する必要がある。

3.4 傾向性の挙動

ここまでが、ロス・エレブモデルの基本的な枠組みである。この枠組みでダイナミクスを考えると、どのようなことが起こるのであろうか。傾向性パラメーター p_{ij} の挙動、戦略採用確率 x_{ij} の挙動の順に調べていこう。

まず、傾向性パラメーターの挙動を大まかに考えてみよう。式2において、強化パラメーター $R_{ij}(t)$ がもし一定の値だとすると、 p_{ij} はどのようになるのであろうか。

このとき

$$p_{ij}(t+1) = (1 - \quad) p_{ij}(t) + R_{ij}$$

$$= p_{ij}(t) + R_{ij} - p_{ij}(t)$$

より

$$p_{ij}(t+1) - p_{ij}(t) = R_{ij} - p_{ij}(t)$$

である。

これより、

$$R_{ij} > p_{ij}(t) \text{ のとき } p_{ij} \text{ は増加}$$

$$R_{ij} < p_{ij}(t) \text{ のとき } p_{ij} \text{ は減少}$$

することがわかる。

p_{ij} が定常となるのは、 $R_{ij} = p_{ij}(t)$ つまり、

$$p_{ij}(t) = R_{ij} / \lambda \quad (\text{式 4})$$

のときである。 p_{ij} がこれより大きいと p_{ij} は減少、小さいと増加するので、式 4 の p_{ij} は漸近安定である。

これより、強化 R_{ij} が一定の場合は傾向パラメーター $p_{ij}(t)$ は強化の値を忘却パラメーター λ で割った値に収束することが分かる。

これが傾向性パラメーターの大ききな挙動である。強化値が大きいほど収束値は大きくなるし、忘却の割合が大きいほど収束値は小さくなる。強化値は実際には一定ではなく、ゲームの繰り返しごとに変化するはずであるが、強化値の平均が大きいと収束値が大きき、小さいと小さきことが予想される。

3.5 採用確率の挙動

次に、プレイヤー i による戦略 j の採用確率 x_{ij} の挙動を考えよう。ここでは、簡単のために戦略は A と B の二つの場合を考える。このとき

$$x_{ia}(t) = p_{ia}(t) / (p_{ia}(t) + p_{ib}(t))$$

である。2 戦略では、 $x_{ia}(t)$ の挙動を考えれば十分なので以下ではこれを単に $x(t)$ と表記する。また、傾向性についても添え字 i を省略して書くことにすると

$$x(t) = p_a(t) / (p_a(t) + p_b(t))$$

となる。

ここで計算上の工夫として

$$k(t) = p_a(t) + p_b(t)$$

と置くことにする。このとき

$$x(t) = p_a(t) / k(t)$$

$$p_a(t) = x(t)k(t)$$

$$p_b(t) = (1 - x(t))k(t)$$

となる。

さて、以下では時刻 t と時刻 $t + 1$ の間の x の変化 Δx を考えることにしよう。

$$\begin{aligned} \Delta x &= x(t+1) - x(t) \\ &= p_a(t+1) / k(t+1) - x(t) \\ &= ((1 - x(t))p_a(t+1) - x(t)p_b(t+1)) / k(t+1) \end{aligned} \quad (\text{式5})$$

である。

ここで、

$$\begin{aligned} p_a(t+1) &= (1 - \lambda) p_a(t) + R_a \\ &= (1 - \lambda) x(t) k(t) + R_a \\ p_b(t+1) &= (1 - \lambda) p_b(t) + R_b \\ &= (1 - \lambda) (1 - x(t)) k(t) + R_b \end{aligned}$$

なので、これを式5の分子に代入すると

$$\begin{aligned} \text{式5の分子} &= (1 - x(t)) \lambda ((1 - \lambda) x(t) k(t) + R_a) \\ &\quad - x(t) \lambda ((1 - \lambda) (1 - x(t)) k(t) + R_b) \\ &= (1 - x(t)) \lambda (1 - \lambda) x(t) k(t) + (1 - x(t)) \lambda R_a \\ &\quad - x(t) \lambda (1 - \lambda) (1 - x(t)) k(t) - x(t) \lambda R_b \\ &= (1 - x(t)) \lambda R_a - x(t) \lambda R_b \end{aligned}$$

となる。

以上より

$$\Delta x = ((1 - x(t)) \lambda R_a - x(t) \lambda R_b) / k(t+1) \quad (\text{式6})$$

となることが分かる。式6を見れば戦略Aの採用確率 $x(t)$ の変化を知ることができる。

3.6 調整ゲームのダイナミクス

式6は具体的にはどのようなダイナミクスなのであろうか。次の調整ゲームを例にとって考えてみよう。

自分 \ 相手	A	B
A	3	1
B	1	2

ここで、利得表の数値はそれぞれの場合に戦略AやBを採ろうとする傾向性に与えられる強化値を表しているものとする。例えば、自分が戦略A、相手が戦略Aを採ったときに自分が戦略Aを採ろうとする傾向性に与えられる強化値 R_a は3である。このとき、戦略Bは採用されていないので、戦略Bを採ろうとする傾向性に与えられる強化値 R_b は0である。

このゲームを自分と相手の二人で繰り返してプレーする場合を考える。自分がAを採る確率を x 、相手がAを採る確率を y とすると、自分と相手の戦略の組み合わせ（戦略プロファイル）は

確率 $x y$ で (A, A)

確率 $x(1 - y)$ で (A, B)

確率 $(1 - x)y$ で (B, A)

確率 $(1 - x)(1 - y)$ で (B, B)

となる。これより R_a 、 R_b の値は

確率 $x y$ で $R_a = 3$ 、 $R_b = 0$

確率 $x(1 - y)$ で $R_a = 1$ 、 $R_b = 0$

確率 $(1 - x)y$ で $R_a = 0$ 、 $R_b = 1$

確率 $(1 - x)(1 - y)$ で $R_a = 0$ 、 $R_b = 2$

である。

したがって、式6より x は

確率 $x y$ で $3(1 - x)/((1 -)k(t) + 3)$

確率 $x(1 - y)$ で $(1 - x)/((1 -)k(t) + 1)$

確率 $(1 - x)y$ で $-x/((1 -)k(t) + 1)$

確率 $(1 - x)(1 - y)$ で $-2x/((1 -)k(t) + 2)$

となる。

分母が少しずつ異なるが、ダイナミクスの概略を知るために大体同じ値

と考えると x の期待値を求めると、次のようになる。ちなみに、 β が十分小さな値のときには、式 4 より $k(t) = p_a(t) + p_b(t)$ が大きな値となるので、分母が同じ値と考えても近似的には差し支えない。

このとき、分母を一定としたときの x の期待値は

$$\begin{aligned} x \text{ の期待値} &= (x y \cdot 3(1-x) + x(1-y) \cdot (1-x) \\ &\quad - (1-x)y \cdot x - (1-x)(1-y) \cdot 2x) / \text{分母} \\ &= x(1-x)(3y + (1-y) - y - 2(1-y)) / \text{分母} \\ &= x(1-x)(3y - 1) / \text{分母} \quad (\text{式7}) \end{aligned}$$

という簡単な式になる。

同様な計算を y についても行なうと

$$y \text{ の期待値} = y(1-y)(3x - 1) / \text{分母} \quad (\text{式8})$$

となる。この二つの式が調整ゲームを自分と相手の二人で繰り返し行なった場合の試行錯誤ダイナミクス方程式となる。

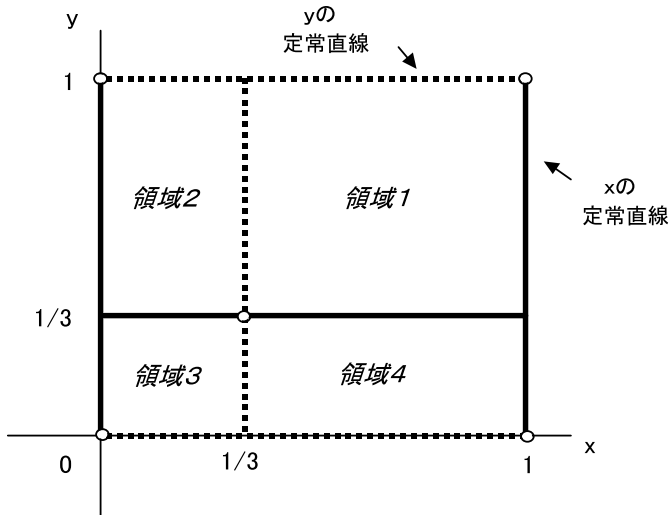


図3 試行錯誤ダイナミクスの定常直線と定常点

式7、式8から、ベクトル図を描いてみよう。このダイナミクスは自分がAを採る確率 x と相手がAを採る確率 y の二つの変数を持つので、状態空間は直線ではなく平面となる。また x と y の範囲は

$$\begin{matrix} 0 & x & 1 \\ 0 & y & 1 \end{matrix}$$

であるので、状態空間は $(0, 0)$ $(0, 1)$ $(1, 0)$ $(1, 1)$ を頂点とする正方形になる。

まず定常点を考えると、式7から x の期待値 = 0 となるのは

$$x = 0, x = 1, y = 1/3$$

である。これを図3の太線で示す。また、式8から y の期待値 = 0 となるのは

$$y = 0, y = 1, x = 1/3$$

である。これを図3の点線で示す。この太線と点線の交点がダイナミクスの定常点であるが、それは状態空間の4頂点と $(1/3, 1/3)$ の5点となる。

次にダイナミクスの方向を矢印で記入する。定常直線で状態空間は四つの領域に分割されるが、それぞれの領域で

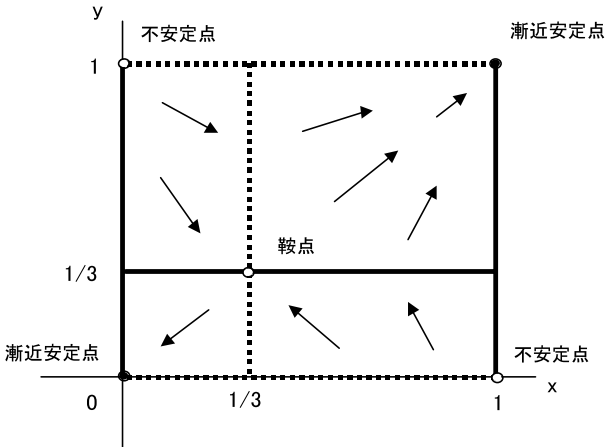


図4 試行錯誤ダイナミクスのベクトル図

	x の期待値	y の期待値
領域 1	+	+
領域 2	+	-
領域 3	-	-
領域 4	-	+

となるので、それぞれ右上、右下、左下、左上に矢印を記入すると図 4 のベクトル図が出来上がる。これが、二人調整ゲームの試行錯誤ダイナミクスを示すベクトル図である。

これより、

$(0, 0)$ $(1, 1)$ が漸近安定点

$(0, 1)$ $(1, 0)$ が不安定点

$(1/3, 1/3)$ が鞍点

となることが分かる。 $(1, 1)$ は両者が A、 $(0, 0)$ は両者が B を採る状態で、調整ゲームの strict Nash 均衡である。したがって、試行錯誤ダイナミクスでも、strict Nash 均衡が漸近安定となることが予想される。

3.7 2 × 2 ゲームのダイナミクス

上の結果は、 x や y の期待値を求めることで得られた結果である。実際には、 x や y は確率的に変化するので期待値の方向にそのまま変化していくわけではなく、変化の方向にはある程度の散らばりがある。しかし、期待値の方向がダイナミクスの収束を示しているときには、変化の方向にある程度の散らばりがあっても、それが大きすぎない限り実際のダイナミクスもやはり収束する。したがって、 x や y の期待値を求めて分析する前節の方法は、ダイナミクスの収束が起きる場合には有用である。

ここでは、前節の方法を用いて一般の 2×2 対称ゲームのダイナミクスを考えてみよう。次のゲームを考える。a、b、c、d はそれぞれの場合の強化を表す正の数である。

自分 \ 相手	A	B
A	a	b
B	c	d

自分がAを採る確率を x 、相手がAを採る確率を y とする。このとき、自分がAやBを採ろうとする傾向パラメーターに与えられる強化値 R_a 、 R_b はそれぞれ

$$\begin{aligned} \text{確率 } x \cdot y \text{ で} & R_a = a, R_b = 0 \\ \text{確率 } x(1 - y) \text{ で} & R_a = b, R_b = 0 \\ \text{確率 } (1 - x)y \text{ で} & R_a = 0, R_b = c \\ \text{確率 } (1 - x)(1 - y) \text{ で} & R_a = 0, R_b = d \end{aligned}$$

となる。

このとき、 x は式6より

$$\begin{aligned} \text{確率 } x \cdot y \text{ で} & a(1 - x) / ((1 - x)k(t) + a) \\ \text{確率 } x(1 - y) \text{ で} & b(1 - x) / ((1 - x)k(t) + b) \\ \text{確率 } (1 - x)y \text{ で} & -cx / ((1 - x)k(t) + c) \\ \text{確率 } (1 - x)(1 - y) \text{ で} & -dx / ((1 - x)k(t) + d) \end{aligned}$$

となる。

ここでも計算の簡略化のために、これらの分母を一定値とみなすと

$$\begin{aligned} x \text{ の期待値} &= (x \cdot y \cdot a(1 - x) + x(1 - y) \cdot b(1 - x) \\ &\quad - (1 - x)y \cdot cx - (1 - x)(1 - y) \cdot dx) / \text{分母} \\ &= x(1 - x)(ay + b(1 - y) - cy - d(1 - y)) / \text{分母} \end{aligned}$$

となる。

ここで $ay + b(1 - y)$ は自分が戦略Aを採ったときの期待効用 u_{1a} に等しい。また、 $cy + d(1 - y)$ は自分が戦略Bを採ったときの期待効用 u_{1b} である。これらを用いると

$$x \text{ の期待値} = x(1 - x)(u_{1a} - u_{1b}) / \text{分母} \quad (\text{式9})$$

同様に y の期待値について計算すると

$$y \text{ の期待値} = y(1 - y)(u_{2a} - u_{2b}) / \text{分母} \quad (\text{式10})$$

となる。ただし、 u_{2a} は相手がAを採ったときの相手の期待効用、 u_{2b} は相手がBを採ったときの相手の期待効用である。

式9、式10が一般の 2×2 対称ゲームのダイナミクス方程式である。これより、試行錯誤ダイナミクスでは

- 1) $(0, 0)$ $(0, 1)$ $(1, 0)$ $(0, 0)$ の各頂点が定常点
- 2) 期待効用の大きい戦略の採用確率が増加することがわかる。
2)は分母が一定とみなした場合の近似的な結果であるが、分かりやすい性質である。

3.8 頂点の安定性

式9、式10から状態空間の頂点が定常であることが分かるが、頂点の安定性についてはどうであろうか。これについては、次の命題が成立する。

【命題3】 二人二戦略ゲームで戦略プロファイル (I, J) がstrict Nash均衡であれば、 (I, J) は試行錯誤ダイナミクスの漸近安定状態である。

【証明】

戦略プロファイル (A, A) がstrict Nash均衡の場合を考えると、このとき $a > c$ である。

$$f_1(x, y) = u_{1a} - u_{1b}$$

$$f_2(x, y) = u_{2a} - u_{2b}$$

を考えると

$$f_1(1, 1) = a - c > 0$$

$$f_2(1, 1) = a - c > 0$$

である。 f_1 、 f_2 は連続関数なので、これより $x = 1$ 、 $y = 1$ の近傍に

$$f_1(x, y) > 0$$

$$f_2(x, y) > 0$$

となる領域が存在する。式9、式10より、この領域では

x の期待値 > 0

y の期待値 > 0

なので $x = 1$ 、 $y = 1$ 、つまり戦略プロファイル (A, A) を採る状態は漸近安定である。 (B, B) が strict Nash 均衡の場合も同様に証明できる。

次に戦略プロファイル (A, B) が strict Nash 均衡の場合を考えると、このとき $c > a$ 、 $b > d$ である。これより

$$f_1(1, 0) = b - d > 0$$

$$f_2(1, 0) = a - c < 0$$

となる。 f_1 、 f_2 は連続関数なので $x = 1$ 、 $y = 0$ の近傍に

$$f_1(x, y) > 0$$

$$f_2(x, y) < 0$$

となる領域が存在する。この領域では

x の期待値 > 0

y の期待値 < 0

なので $x = 1$ 、 $y = 0$ 、つまり戦略プロファイル (A, B) を採る状態は漸近安定である。 (B, A) が strict Nash 均衡の場合も同様に証明できる。

【証明終わり】

この命題は、頂点に対応する戦略プロファイルが strict Nash 均衡であれば、その頂点が漸近安定であることを示している。逆に、頂点に対応する戦略プロファイルが Nash 均衡でなければ、次の命題が成立する。

【命題 4】 二人二戦略ゲームで戦略プロファイル (I, J) が Nash 均衡でなければ、 (I, J) は試行錯誤ダイナミクスで不安定である。

【証明略】

これらの命題は、静学解を吟味することで試行錯誤ダイナミクスにおける安定性を判定できる場合があることを示している。

3.9 二人モデルと集団モデル

3.6から3.9までで考察してきたモデルは二人のプレーヤーが繰り返し何回かゲームを行なうことを想定したモデルである。これに対し、レプリケーターダイナミクスや模倣ダイナミクスで想定してきたのは、プレーヤーの集団があってその中からプレーヤーがランダムに選ばれてゲームを行なう状況であった。

後者では集団中の戦略のシェア（頻度）を考え、シェアの変化を考えるとというアプローチが可能であったのに対し、前者ではプレーヤーが二人しかいないため「戦略のシェア」を考えるアプローチは余り意味がない。そこで、前者では「戦略のシェア」の代わりに各々のプレーヤーがある戦略を採用する確率を考え、この「採用確率」のダイナミクスを考える、というアプローチを取った。

このように、両者はダイナミクスを考えるという点では共通しているが、想定する状況や、変化するものの中味が大きく異なっているため、異なるタイプのモデルと考える方がよい。ここでは、前者のタイプのモデルを二人モデル（一般には非集団モデル）、後者のタイプのモデルを集団モデルと呼ぶことにしよう。

社会の中である考え方や行動のパターンが広まっていったり消滅したりする現象をモデル化するには、集団モデルが適している。一方、夫婦の間で役割分担が出来上がっていくプロセスや、寡占企業同士の駆け引き、隣国同士の相互作用などを考えるときには、ランダムマッチングモデルを用いるよりもここで紹介した二人モデルの方が適している。

このように、いずれのモデルにも分析に適した社会現象がある。試行錯誤ダイナミクスについてはここまで二人モデルのみを紹介してきたが、では集団モデル型の試行錯誤モデルは考えられないのであろうか。以下の節ではこの点について考えてみよう。

3.10 集団型試行錯誤ダイナミクス

2 × 2 ゲーム

自分 \ 相手	A	B
A	a	b
B	c	d

をランダムマッチングで行なっているプレーヤーの集団を考える。各プレーヤーは試行錯誤で戦略を修正しているものとする。このとき、集団内の戦略採用頻度はどのように変化するのであろうか。

プレーヤー i が戦略 j を採ろうとする傾向性 $p_{ij}(t)$ が

$$p_{ij}(t+1) = (1 - \alpha) p_{ij}(t) + \alpha R_{ij}(t)$$

にしたがって変化しているとする。このとき、プレーヤー i が戦略 A を採る確率を $x_i(t)$ とすると、これの変化 \dot{x}_i は

$$\dot{x}_i = (\alpha (1 - x_i(t)) R_{ia} - \alpha x_i(t) R_{ib}) / (p_{ia}(t+1) + p_{ib}(t+1))$$

となる。ここまでは 3.5 節の考察と同じである。

ここで、集団中の x_i の平均として x を考えると、集団中からランダムに一人選んだ相手が戦略 A を採る確率は x となる。このとき、プレーヤー i が A や B を採ろうとする傾向パラメーターに与えられる強化値 R_{ia} 、 R_{ib} はそれぞれ

$$\text{確率 } x_i x \text{ で} \quad R_{ia} = a, R_{ib} = 0$$

$$\text{確率 } x_i (1 - x) \text{ で} \quad R_{ia} = b, R_{ib} = 0$$

$$\text{確率 } (1 - x_i) x \text{ で} \quad R_{ia} = 0, R_{ib} = c$$

$$\text{確率 } (1 - x_i) (1 - x) \text{ で} \quad R_{ia} = 0, R_{ib} = d$$

となる。

これより、 \dot{x}_i は

$$\text{確率 } x_i x \text{ で} \quad (1 - x_i) a / ((1 - \alpha) k_{ia}(t) + a)$$

$$\text{確率 } x_i (1 - x) \text{ で} \quad (1 - x_i) b / ((1 - \alpha) k_{ia}(t) + b)$$

$$\text{確率 } (1 - x_i) x \text{ で} \quad -x_i c / ((1 - \alpha) k_{ib}(t) + c)$$

確率 $(1 - x_i)(1 - x)$ で $-x_i d / ((1 - x)k_i(t) + d)$

となる。

分母が少しずつ異なるが、ここでもこれらをおおむね同じと考えて x_i の期待値を求めると

$$\begin{aligned} x_i \text{の期待値の分子} &= x_i x (1 - x_i) a + x_i (1 - x)(1 - x_i) b \\ &\quad - (1 - x_i) x x_i c - (1 - x_i)(1 - x) x_i d \\ &= x_i (1 - x_i) [x a + (1 - x) b - x c - (1 - x) d] \end{aligned}$$

ここで $x a + (1 - x) b$ はプレイヤー i が戦略 A を採ったときの期待効用 u_a に等しく、 $x c + (1 - x) d$ は戦略 B を採ったときの期待効用 u_b に等しいので、

$$x_i \text{の期待値} = x_i (1 - x_i) (u_a - u_b) / \text{分母} \quad (\text{式11})$$

となる。

これより $u_a > u_b$ のとき x_i の期待値は増加し、 $u_a < u_b$ のとき x_i の期待値は減少することがわかる。ところで、この結果は任意のプレイヤー i について成立する。ここで

$$x = x_i / N$$

であるので、任意のプレイヤーについて $x_i > 0$ ならば $x > 0$ 、 $x_i < 0$ ならば $x < 0$ である。

以上より、

$$u_a > u_b \text{ のとき } x > 0$$

$$u_a < u_b \text{ のとき } x < 0$$

であることがわかる。

3.11 集団型試行錯誤ダイナミクスの性質

式11から、集団型試行錯誤ダイナミクスの性質について幾つかの命題が導ける。まず、状態空間の頂点の定常性についてである。ちなみに、2戦略集団モデルでは、集団の中で戦略 A が採用される確率の平均 x が分かれば、戦略 B が採用される確率の平均 $1 - x$ も分かるので、状態空間は $0 \leq x \leq 1$ の線分(単位単体)となる。状態空間の頂点は、単位単体の両端 $x = 0$ と

$x = 1$ である。

状態空間の頂点について次の命題が成立する。

【命題 5】 2×2 対称ゲームの集団型試行錯誤ダイナミクスにおいて、状態空間の頂点は定常点である。

【証明】

$x = 0$ の場合を考えると、 $x = x_i / N$ と $x_i = 0$ より、 $x = 0$ のときは任意のプレイヤー i について $x_i = 0$ である。このとき式11より任意の i について $x_i = 0$ である。したがって $x = 0$ 。ゆえに $x = 0$ は定常点である。

$x = 1$ の場合は、 $0 < x_i < 1$ より任意の i について $x_i = 1$ である。したがって、任意の i について $x_i = 0$ なので $x = 0$ 。ゆえに $x = 1$ は定常点である。

【証明終わり】

二人モデルの場合と同様に、集団モデルの場合も状態空間の頂点は定常点である。頂点の安定性については次の命題が成立する。

【命題 6】 2×2 対称ゲームにおいて、戦略プロファイル (I, I) が strict Nash 均衡であるならば、戦略 I に対応する頂点は、集団型試行錯誤ダイナミクスにおいて漸近安定である。

【証明】

戦略プロファイル (A, A) が strict Nash 均衡の場合を考える。このとき $a > c$ である。式11において $f(x) = u_a - u_b$ とおくと

$$f(1) = a - c > 0$$

である。

$f(x)$ は連続な関数なので、 $f(1) > 0$ より $x = 1$ の近傍に

$$f(x) > 0$$

となる領域があることがわかる。したがって、 x が 1 から少しずれても、 x が十分小さければ

$$x_i \text{ の期待値} > 0$$

となる。ゆえに、戦略 A に対応する頂点 $x = 1$ は漸近安定である。

戦略プロファイル (B, B) が strict Nash 均衡の場合も同様に証明できる。

【証明終わり】

【命題 7】 2×2 対称ゲームにおいて、戦略プロファイル (I, I) が Nash 均衡でなければ、戦略 I に対応する頂点は、集団型試行錯誤ダイナミクスにおいて不安定である。

【証明略】

例えば、戦略プロファイル (A, A) が Nash 均衡でなければ、 $a < c$ となることから証明できる。

このように、各プレーヤーが試行錯誤で戦略を修正する場合の平均戦略採用頻度に関しても、レプリケーターダイナミクスや模倣ダイナミクスと同様の命題が成立する。

4 まとめ

本論文では学習ダイナミクスのうち、模倣ダイナミクスと試行錯誤ダイナミクスについて検討してきた。これらはいずれも、プレーヤーが利得構造について知らない場合にも実現可能なダイナミクスである。 2×2 ゲームについての分析の結果、いずれのダイナミクスにおいても strict Nash 均衡に対応する状態が漸近安定となり、Nash 均衡ではない状態が不安定となることが示された。

模倣ダイナミクスと試行錯誤ダイナミクスは現実の社会的相互作用においても機能している可能性が高いダイナミクスなので、これらのダイナミ

クスモデルの整備によって、より広い範囲の社会現象がゲーム理論を用いて分析できるようになることが期待される。

文献

Erev, I. & Roth, A. E. 1998. "Predicting How People Play Games: Reinforcement Learning in Experimental Games with Unique, Mixed Strategy Equilibria" *The American Economic Review* 88 (4) 848-881

石原英樹, 金井雅之 2002. 『進化的意思決定』 朝倉書店

Roth, A. E. & Erev, I. 1995. "Learning in Extensive-Form Games: Experimental Data and Simple Dynamic Models in the Intermediate Term" *Games and Economic Behavior* 8 164-212