

報酬駆動型システムにおける報酬の設計と報酬による最適化

奥原 浩之*

1. はじめに

個体の最適な応答の結果として、集団にとって効率的な選択が実現されることは望ましい。しかしながら、環境によっては、個体の合理的な意思決定で集団にとって望ましい状況が実現されるとは限らない。そこで、環境を人為的に設計することで、設計者が望ましいと考える集団として効率的な選択を、個体の最適な応答により自律分散的に実現する仕組みが求められる。

このような背景から、特定の戦略の選択を促進・抑制するメカニズムを、ペイオフの観点から分析することが求められる[1]。遺伝子、神経回路、進化ゲームに関する力学[1–4]は、意思決定の主体が戦略を選び行動した結果、報酬を原動力として駆動していると見なすことができ、ペイオフにもとづく効用関数の変化を報酬として入力している報酬駆動型システムと捉えることができる。

本解説では、その中でも、まず効用関数の最大化が適者生存則など[5,6]に関連していることを紹介する。つぎに、与えられたペイオフのもとで均衡が漸近安定とならない場合に、ペイオフを変更することで漸近安定化[7,8]する報酬の設計法を概説する。そして、逐次得られる限られたペイオフのもとで、最もペイオフが高い戦略を探索する報酬による最適化[9,10]について概観する。

2. メカニズムデザインとペイオフ

2.1 非協力ゲームとペイオフ双行列

本解説で想定しているメカニズムの設計について、二人非協力有限ゲームにおける標準型ゲームを例にして簡単に述べる。

意思決定の主体であるプレイヤーの集合が $P = \{A, B\}$ 、プレイヤー $i (\in P)$ の戦略の集合を $S_i = \{s_i^1, s_i^2\}$ とし、選択した行動を $s_i (\in S_i)$ とする。プレイヤー i が戦略 π を選択することを s_i^π で表す。プレイヤー A, B がゲームをプレーすることで得られる利得であるペイオフ双行列 $G = (V_{AB}, V_{BA})$ を第1表に示す。

ここで、

$$V_{AB} = \begin{pmatrix} v_{AB}^{11} & v_{AB}^{12} \\ v_{AB}^{21} & v_{AB}^{22} \end{pmatrix}, \quad V_{BA} = \begin{pmatrix} v_{BA}^{11} & v_{BA}^{21} \\ v_{BA}^{12} & v_{BA}^{22} \end{pmatrix} \quad (1)$$

第1表 ペイオフ双行列の例

	s_B^1	s_B^2
s_A^1	$(v_{AB}^{11}, v_{BA}^{11})$	$(v_{AB}^{12}, v_{BA}^{21})$
s_A^2	$(v_{AB}^{21}, v_{BA}^{12})$	$(v_{AB}^{22}, v_{BA}^{22})$

である。

ペイオフの値によって、囚人のジレンマ、鷹鳩ゲームなどの異なる状況を表現することになる。

プレイヤーが N 人となる場合を考える。このとき、プレイヤーの選択した行動 $s^* = [s_1^* s_2^* s_3^* \cdots s_N^*]^T \in \mathbb{R}^{N \times 1}$ が互いに最適な応答

$$s_i^* \in \arg \max_{s_i \in S_i} u_i(s_i, s_i^*) \quad (\forall i) \quad (2)$$

となるとき、Nash 均衡であるという。ここで、 s_i^* は s^* から s_i^* を除いた $[s_1^* \cdots s_{i-1}^* s_{i+1}^* \cdots s_N^*]^T \in \mathbb{R}^{(N-1) \times 1}$ であり、効用関数 u_i は戦略空間 $S = \prod_{i \in N} S_i$ からペイオフへの写像を与える。

たとえば、第1表の場合、よく知られているように、 $V_{AB} = V_{BA}^T$ で、 $v_{AB}^{21} > v_{AB}^{11} > v_{AB}^{22} > v_{AB}^{12}$ かつ $2v_{AB}^{11} > v_{AB}^{21} + v_{AB}^{12}$ であるなら、Nash 均衡は (s_A^1, s_B^2) となる。このことは、個人の最適な応答が集団にとって効率的な選択 (s_A^1, s_B^1) とならない囚人のジレンマの状況を表している。

2.2 混合戦略と複製方程式

プレイヤー i が戦略 π を頻度 x_i^π で選択する混合戦略の場合を考え、 $\mathbf{x}_i = [x_i^1 x_i^2 x_i^3 \cdots x_i^{n_i}]^T \in \mathbb{R}^{n_i \times 1}$ とする。ここで、 n_i はプレイヤー i の戦略の数であり $\sum_{\pi=1}^{n_i} x_i^\pi = 1$ である。このとき、純戦略 s_i^π の選択は $x_i^\pi = 1$ であるので、 \mathbf{x}_i において $x_i^\pi = 1$ かつ $x_i^{\pi'} = 0$ ($\pi' \neq \pi$)とした $\mathbf{e}_i^\pi = [0 \cdots 0 1 \cdots 0]^T \in \mathbb{R}^{N \times 1}$ と表せる。

プレイヤー i とそれ以外のプレイヤー $j (\neq i)$ が、ゲームをプレーするときの効用関数 u_i が期待利得で与えられるとすると、

$$u_i(\mathbf{x}_i, \mathbf{x}_{\bar{i}}) = \sum_{j=1, j \neq i}^N \mathbf{x}_i^T V_{ij} \mathbf{x}_j \quad (3)$$

* 大阪大学大学院 情報科学研究科

Key Words: reward-driven system, design, optimization.

となる。ここで、 $\mathbf{x}_i = [\mathbf{x}_1^T \cdots \mathbf{x}_{i-1}^T \mathbf{x}_{i+1}^T \cdots \mathbf{x}_N^T]^T \in \Re^{(M-n_i) \times 1}$, $\mathbf{V}_{ij} \in \Re^{n_i \times n_j}$ であり、 $M = \sum_{i=1}^N n_i$ である。

このとき、任意の有限標準形ゲームには混合戦略の範疇で

$$u_i(\mathbf{x}_i^*, \mathbf{x}_{\bar{i}}^*) \geq u_i(\mathbf{e}_i^\pi, \mathbf{x}_{\bar{i}}^*) \quad (\forall i, \forall \pi) \quad (4)$$

を必要十分条件とする Nash 均衡 $\mathbf{x}^* = [\mathbf{x}_1^{*\text{T}} \mathbf{x}_2^{*\text{T}} \mathbf{x}_3^{*\text{T}} \cdots \mathbf{x}_N^{*\text{T}}]^T \in \Re^{M \times 1}$ が存在するが、必ずしも唯一とは限らない。

いま、戦略 π の選択による利得が期待利得 $u_i(\mathbf{x}_i, \mathbf{x}_{\bar{i}})$ より高いようであれば、プレイヤー i は戦略の頻度 x_i^π を増加させ、低いようであれば減少させると考える。あるいは、期待利得に比べて利得が高い戦略を選択した集団の割合が複製により増加して、低い戦略を選択した集団の割合が淘汰により減少すると考える。

ゲームの動学的な状況を表す進化ゲームでは、複製方程式はリプリケータダイナミクスとして、

$$\dot{x}_i^\pi = x_i^\pi (u_i(\mathbf{e}_i^\pi, \mathbf{x}_{\bar{i}}) - u_i(\mathbf{x}_i, \mathbf{x}_{\bar{i}})) \quad (\forall i, \forall \pi) \quad (5)$$

とできる[1]。

たとえば、第1表の場合、プレイヤー A の混合戦略は $(x_A^1, x_A^2) = (0, 1)$ 、プレイヤー B の混合戦略は $(x_B^1, x_B^2) = (0, 1)$ となり Nash 均衡である純戦略の組 (s_A^2, s_B^2) に収束することが示される。

3. 報酬駆動型システムの概要

3.1 遺伝子・神経回路網のモデル

ここで、ある種の遺伝子、結合振動子、神経回路網のモデルと勾配系やゲームのモデルにおけるアナロジーについて述べておく。

まず、遺伝子 A^π の相対頻度 x^π が従う連続時間の淘汰・突然変異のモデルの一つに

$$\begin{aligned} \dot{x}^\pi &= x^\pi (u^\pi(\mathbf{e}^\pi, \mathbf{x}) - u(\mathbf{x}, \mathbf{x})) \\ &\quad + \sum_{\pi'=1}^n (\epsilon^{\pi\pi'} x^{\pi'} - \epsilon^{\pi'\pi} x^\pi) \quad (\forall \pi) \end{aligned} \quad (6)$$

がある。ここで、 $\mathbf{x} = [x^1 x^2 x^3 \cdots x^n]^T \in \Re^{n \times 1}$, $u(\mathbf{x}, \mathbf{x}) = \mathbf{x}^T \mathbf{M} \mathbf{x}$ は集団平均適応度、行列 $\mathbf{M} \in \Re^{n \times n}$ の $\pi\pi'$ 要素はマルサス的適応度である。 $\mathbf{e}^\pi = [0 \cdots 010 \cdots 0]^T \in \Re^{n \times 1}$ は \mathbf{x} において $x^\pi = 1$ かつ $x^{\pi'} = 0$ ($\pi' \neq \pi$) としたものである。

$\epsilon^{\pi\pi'}$ は突然変異率であり、 $\epsilon^{\pi\pi'} \geq 0$, $\sum_{\pi=1}^n \epsilon^{\pi\pi'} = 1$ ($\forall \pi'$) を満たしている。右辺第一項が淘汰、第二項が突然変異を表す[1]。

つぎに、非線形振動する素子を結合した振動子が従う連続時間の結合振動子のモデルの一つに

$$\dot{\mathbf{x}}_i = \mathbf{f}(\mathbf{x}_i) + D \sum_{j=1}^N (\mathbf{x}_j - \mathbf{x}_i) \quad (\forall i) \quad (7)$$

がある。右辺第一項は非線形項、第二項は拡散係数 D をもつ線形項である。微分方程式 $\dot{\mathbf{x}}_i = \mathbf{f}(\mathbf{x}_i)$ が写像関数 $\mathbf{x}_i^{(k+1)} = \mathbf{g}(\mathbf{x}_i^{(k)})$ をもつと仮定する。このとき、サンプリング間隔を T ($t = kT$), $\hat{D} = \frac{1 - \exp(-NDT)}{N}$ として、離散時間のモデル

$$\mathbf{x}_i^{(k+1)} = \mathbf{g}(\mathbf{x}_i^{(k)}) + \hat{D} \sum_{j=1}^N (\mathbf{g}(\mathbf{x}_j^{(k)}) - \mathbf{g}(\mathbf{x}_i^{(k)})) \quad (8)$$

が得られる[2]。

離散時間の結合振動子のモデルの両辺に非線形変換を施すと

$$\mathbf{g}(\mathbf{x}_i^{(k+1)}) = \mathbf{g}\left((1-\epsilon)\mathbf{g}(\mathbf{x}_i^{(k)}) + \frac{\epsilon}{N} \sum_{j=1}^N \mathbf{g}(\mathbf{x}_j^{(k)})\right) \quad (9)$$

となる。ここで、 $n_i = n$ ($\forall i$) であり、 $\mathbf{V}_{ii} = \left(1 - \epsilon + \frac{\epsilon}{N}\right) \mathbf{I}_n$, $\mathbf{V}_{ij} = \frac{\epsilon}{N} \mathbf{I}_n$ ($j \neq i$) として、変数変換 $\mathbf{y}_i^{(k)} = \mathbf{g}(\mathbf{x}_i^{(k)})$ を施すと、神経回路網のモデル

$$\begin{aligned} \mathbf{y}_i^{(k+1)} &= \mathbf{g}\left(\sum_{j=1}^N \mathbf{V}_{ij} \mathbf{y}_j^{(k)}\right) \\ &= \mathbf{g}\left(\frac{\partial \tilde{u}_i(\mathbf{y}_i^{(k)}, \mathbf{y}_i^{(k)})}{\partial y_i^{\pi(k)}}\right) \end{aligned} \quad (10)$$

の一例とも見なせる[3]。ただし、

$$\tilde{u}_i(\mathbf{x}_i, \mathbf{x}_{\bar{i}}) = u_i(\mathbf{x}_i, \mathbf{x}_{\bar{i}}) + \frac{1}{2} \mathbf{x}_i^T \mathbf{V}_{ii} \mathbf{x}_i \quad (11)$$

であり、 $\mathbf{I}_n \in \Re^{n \times n}$ は単位行列である。

3.2 勾配系・ N 人進化ゲームのモデル

ところで、プレイヤー i の状態が n_i 個の変数をもつベクトル \mathbf{x}_i で表される N 人の勾配系モデルの一つに

$$\dot{x}_i^\pi = g(x_i^\pi) \phi_i^\pi \quad (\forall i, \forall \pi) \quad (12)$$

がある。ここで、

$$\begin{aligned} \phi_i^\pi &= \sum_{j=1}^N \sum_{\pi'=1}^{n_j} v_{ij}^{\pi\pi'} x_j^{\pi'} + \theta_i^\pi \\ &= \frac{\partial \hat{u}_i(\mathbf{x}_i, \mathbf{x}_{\bar{i}})}{\partial x_i^\pi} \end{aligned} \quad (13)$$

であり、

$$\dot{u}_i(\mathbf{x}_i, \mathbf{x}_{\bar{i}}) = \tilde{u}_i(\mathbf{x}_i, \mathbf{x}_{\bar{i}}) + \theta_i^\pi \mathbf{x}_i \quad (14)$$

である。ただし、 $\theta_i = [\theta_i^1 \theta_i^2 \theta_i^3 \cdots \theta_i^{n_i}]^T \in \Re^{n_i \times 1}$ であり、 $g(x_i^\pi) > 0$ なら値 x_i^π の更新は効用関数 $\hat{u}_i(\mathbf{x}_i, \mathbf{x}_{\bar{i}})$ の値を増加させる。もちろん、均衡が漸近安定であるかどうかは、定常点近傍でのヤコビ行列の固有値により安定判別がなされる[4]。

【定理 1】 効用最大化と適者生存との関係

全プレイヤーの効用関数 $\hat{u}_i(\mathbf{x}_i, \mathbf{x}_{\bar{i}})$ の和の最大化は、任意のベクトル $\mathbf{y}_k \in \mathbb{R}^{m_k \times 1}$, ($k = 1, 2, 3, \dots, S$) に対する供給関数 $\eta(\mathbf{y}_k)$ と、パラメータ φ_i^π をもつ利用関数 $\xi(\mathbf{y}_k, \varphi_i^\pi)$ が与えられるとき、競争係数を

$$\begin{aligned}\lambda_{ij}^{\pi\pi'} &= \lambda_{ji}^{\pi'\pi} = \frac{v_{ij}^{\pi\pi'} + v_{ji}^{\pi'\pi}}{2} \\ &= -\sum_{k=1}^S \xi(\mathbf{y}_k, \varphi_i^\pi) \xi(\mathbf{y}_k, \varphi_j^{\pi'})\end{aligned}\quad (15)$$

とし、適応係数を

$$\theta_i^\pi = \sum_{k=1}^S \eta(\mathbf{y}_k) \xi(\mathbf{y}_k, \varphi_i^\pi) \quad (16)$$

として、供給関数の利用関数による関数近似

$$\eta(\mathbf{y}_k) = \sum_{i=1}^N \sum_{\pi=1}^{n_i} x_i^\pi \xi(\mathbf{y}_k, \varphi_i^\pi) \quad (17)$$

と等価である。

(証明)

$$\begin{aligned}&-\frac{1}{2} \sum_{k=1}^S \left(\eta(\mathbf{y}_k) - \sum_{i=1}^N \sum_{\pi=1}^{n_i} x_i^\pi \xi(\mathbf{y}_k, \varphi_i^\pi) \right)^2 \\ &= -\frac{1}{2} \sum_{k=1}^S \eta^2(\mathbf{y}_k) + \sum_{i=1}^N \boldsymbol{\theta}_i^T \mathbf{x}_i - \frac{1}{2} \sum_{i=1}^N \mathbf{x}_i^T \boldsymbol{\Lambda}_{ij} \mathbf{x}_j \\ &= \sum_{i=1}^N \hat{u}_i(\mathbf{x}_i, \mathbf{x}_{\bar{i}}) - \frac{1}{2} \sum_{k=1}^S \eta^2(\mathbf{y}_k)\end{aligned}\quad (18)$$

ここで、 $\boldsymbol{\Lambda}_{ij} \in \mathbb{R}^{n_i \times n_j}$ は要素に $\lambda_{ij}^{\pi\pi'}$ をもつ。

関数近似の観点から、ベクトル x_i^π やパラメータ φ_i^π に対する適者生存型の更新則 [5] や、プレイヤー数や要素数が変化する動的な環境でも適用可能な自由エネルギー型の更新則 [6] が導かれている。□

さらに、勾配系モデルを N 人進化ゲームとするために、混合戦略の条件 $\sum_{\pi=1}^{n_i} x_i^\pi = 1$ ($\forall i$) を満たすダイナミクスを導出すると

$$\begin{aligned}\dot{x}_i^\pi &= x_i^\pi (\hat{u}_i(\mathbf{e}_i^\pi, \mathbf{x}_{\bar{i}}) - \hat{u}_i(\mathbf{x}_i, \mathbf{x}_{\bar{i}})) \\ &= x_i^\pi \left(\phi_i^\pi - \sum_{\pi'=1}^{n_i} x_i^{\pi'} \phi_i^{\pi'} \right) \\ &= \sum_{\pi'=1}^{n_i} g^{\pi'}(x_i^\pi) \phi_i^{\pi'}\end{aligned}\quad (\forall i, \forall \pi) \quad (19)$$

となる。ここで、 $\boldsymbol{\phi}_i = [\phi_i^1 \phi_i^2 \phi_i^3 \cdots \phi_i^{n_i}]^T \in \mathbb{R}^{n_i \times 1}$ であり、

$$g^{\pi'}(x_i^\pi) = \begin{cases} x_i^\pi (1 - x_i^\pi) & (\pi' = \pi) \\ -x_i^\pi x_i^{\pi'} & (\pi' \neq \pi) \end{cases} \quad (20)$$

である。このとき、初期値が混合戦略の条件を満たしていれば、 $\sum_{\pi=1}^{n_i} \dot{x}_i^\pi = 0$ ($\forall i$) となる。

以上に述べた遺伝子、神経回路網や進化ゲームのモデルには、複製方程式から派生したモデルも含まれており、これらは共通して、戦略を選び行動した結果、ペイオフにもとづく効用関数の変化を入力として動作していると考えることもできる。本解説では、それらを報酬駆動型システムとして考えている。

4. 報酬の設計

4.1 非線形大域結合システム

ここで、遺伝子、神経回路網や進化ゲームの連続時間のモデルである式 (6), (7), (19) をまとめて取り扱うことができる非線形大域結合システムによるモデル化と漸近安定化に向けた制御、報酬の設計について議論する。

システムは N 個の局所システムから構成され、 i 番目の局所システムの状態は n_i 個の変数をもつベクトル \mathbf{x}_i で表される。そのとき、局所システムの状態方程式を

$$\dot{\mathbf{x}}_i = \sum_{j=1}^N \mathbf{G}_{ij}(\mathbf{x}) \{ \mathbf{f}(\mathbf{x}_j) + \sum_{k=1}^N c_{jk} a_{jk} (\boldsymbol{\Gamma}_{jk} \mathbf{x}_k - \boldsymbol{\Gamma}_{jj} \mathbf{x}_j) \} \quad (\forall i) \quad (21)$$

とする。ここで、 $\mathbf{G}_{ij}(\mathbf{x}) \in \mathbb{R}^{n_i \times n_j}$, $\mathbf{f}(\mathbf{x}_j)$ は C^1 級で、 $\mathbb{R}^{n_j} \rightarrow \mathbb{R}^{n_j}$ の非線形写像であり、 $\mathbf{f}(\mathbf{0}_{n_j}) = \mathbf{0}_{n_j}$ である。 $\mathbf{x} = [\mathbf{x}_1^T \mathbf{x}_2^T \cdots \mathbf{x}_N^T]^T \in \mathbb{R}^{M \times 1}$ であり、 $\mathbf{0}_{n_j} \in \mathbb{R}^{n_j \times 1}$ は零ベクトルである。また、 $c_{jk} \in \mathbb{R}^{N \times N}$ は局所システム間の結合強度、 $a_{jk} \in \mathbb{R}^{N \times N}$ は局所システム間の結合情報であり、 $a_{jk} = 1$ は結合あり、 $a_{jk} = 0$ は結合なしを意味する。 $\boldsymbol{\Gamma}_{jk} \in \mathbb{R}^{n_j \times n_k}$ は局所システム変数間の相互作用である。

$$\begin{aligned}\text{いま}, \tilde{a}_{jk} &= a_{jk} \ (k \neq j), \ \tilde{a}_{jj} = -\sum_{k=1}^N a_{jk}, \ \tilde{c}_{jk} = c_{jk} \\ (k \neq j), \ \tilde{c}_{jj} &= \frac{1}{\tilde{a}_{jj}} \sum_{k=1}^N c_{jk} a_{jk} \text{ とすると}\end{aligned}$$

$$\dot{\mathbf{x}}_i = \sum_{j=1}^N \mathbf{G}_{ij}(\mathbf{x}) \tilde{\mathbf{f}}(\mathbf{x}_j) \quad (\forall i) \quad (22)$$

となる。ここで、

$$\tilde{\mathbf{f}}(\mathbf{x}_j) = \mathbf{f}(\mathbf{x}_j) + \sum_{k=1}^N \tilde{c}_{jk} \tilde{a}_{jk} \boldsymbol{\Gamma}_{jk} \mathbf{x}_k \quad (23)$$

である。

システム全体の状態方程式は

$$\dot{\mathbf{x}} = \mathbf{G}(\mathbf{x}) \{ \mathbf{f}(\mathbf{x}) + \boldsymbol{\Gamma}_\alpha \mathbf{x} \} \quad (24)$$

となる。ここで、 $\mathbf{G}(\mathbf{x}) \in \mathbb{R}^{M \times M}$ は ij ブロックに $\mathbf{G}_{ij}(\mathbf{x}) \in \mathbb{R}^{n_i \times n_j}$, $\boldsymbol{\Gamma}_\alpha \in \mathbb{R}^{M \times M}$ は ij ブロックに $\alpha_{ij} \boldsymbol{\Gamma}_{ij} \in \mathbb{R}^{n_i \times n_j}$

をもち、 $\alpha_{ij} = \tilde{c}_{ij}\tilde{a}_{ij}$ である。また、 $\mathbf{f}(\mathbf{x}) = [\mathbf{f}(\mathbf{x}_1)^T \mathbf{f}(\mathbf{x}_2)^T \mathbf{f}(\mathbf{x}_3)^T \cdots \mathbf{f}(\mathbf{x}_N)^T]^T \in \mathbb{R}^{M \times 1}$ である。

4.2 漸近安定化

非線形大域結合システムの漸近安定化は、非線形関数の線形緩和とフィードバック制御により行われる。まず、局所システムの変数 \mathbf{x}_i の均衡 \mathbf{x}_i^* に対して領域

$$\mathbf{D} = \bigcup_{i=1}^N \mathbf{D}_i \quad \mathbf{D}_i = \{\mathbf{x}_i : \|\mathbf{x}_i - \mathbf{x}_i^*\| < \alpha, \alpha > 0\} \quad (25)$$

を考える。

各局所システムに正定対称行列 $\mathbf{P}_i \in \mathbb{R}^{n_i \times n_i}$ をもつ Lyapunov 関数 $V_i(\mathbf{x}_i) = \mathbf{x}_i^T \mathbf{P}_i \mathbf{x}_i$ の存在を仮定し、Lyapunov 関数 $V(\mathbf{x}) = \sum_{i=1}^N V_i(\mathbf{x}_i) : \mathbf{D} \subseteq \mathbb{R}^M \rightarrow \mathbb{R}_+^M$ は C^1 級で $V(\mathbf{x}^*) = 0$, $\mathbf{x}^* \in \mathbf{D}$ とする。

いま、非線形関数の線形緩和のために

$$\{\mathcal{L}_{\mathbf{h}_\mu} V\}(\mathbf{x}) < 0 \quad (26)$$

を満たす Passivity degree $\mu = [\mu_1 \mu_2 \mu_3 \cdots \mu_N]^T \in \mathbb{R}^{N \times 1}$ を考える [7]。ただし、Lie 微分作用素は

$$\begin{aligned} \{\mathcal{L}_{\mathbf{h}_\mu} V\}(\mathbf{x}) &= \sum_{i=1}^N \frac{\partial V(\mathbf{x})}{\partial \mathbf{x}_i} \mathbf{h}_{\mu_i}(\mathbf{x}_i) \\ &= \nabla_{\mathbf{x}} V^T \mathbf{h}_\mu(\mathbf{x}) \end{aligned} \quad (27)$$

であり、 $\mathbf{x} \in \mathbf{D}$, $\mathbf{x} \neq \mathbf{x}^*$ に対して

$$\mathbf{h}_\mu(\mathbf{x}) = \mathbf{G}(\mathbf{x}) \{ \mathbf{f}(\mathbf{x}) - \boldsymbol{\Gamma}_\mu(\mathbf{x}^* - \mathbf{x}) \} \quad (28)$$

である。ここで、 $\boldsymbol{\Gamma}_\mu \in \mathbb{R}^{M \times M}$ は ii ブロックに $\mu_i \boldsymbol{\Gamma}_{ii} \in \mathbb{R}^{n_i \times n_i}$ をもつブロック対角行列である。

Passivity degree μ による線形緩和の後、限られた局所システム ($i = 1, 2, 3, \dots, m$) に対して、フィードバックゲイン $d_i (> 0)$ を用いる Pinning 制御を行う [8]。

$$\tilde{\mathbf{f}}(\mathbf{x}_j) = \mathbf{f}(\mathbf{x}_j) + \sum_{k=1}^N \alpha_{jk} \boldsymbol{\Gamma}_{jk} \mathbf{x}_k + \mathbf{u}_i \quad (\forall i) \quad (29)$$

ただし、

$$\mathbf{u}_i = -(k_i - 1)d_i \boldsymbol{\Gamma}_{jk} \mathbf{x}_i \quad (i = 1, 2, 3, \dots, m) \quad (30)$$

$$\mathbf{u}_i = 0 \quad (i = m+1, m+2, m+3, \dots, N) \quad (31)$$

である。ここでは、局所システムの番号 i は状況に応じて振り直されていることに注意が必要である。

【補題 1】 非線形大域結合システムの漸近安定化

正定対称行列 $\mathbf{P} \in \mathbb{R}^{M \times M}$ をもつ Lyapunov 関数 $V(\mathbf{x}) = \mathbf{x}^T \mathbf{P} \mathbf{x}$ と Passivity degree μ の存在を仮定する。このとき、Pinning 制御を行い

$$\mathbf{M}(\mathbf{x}) \prec 0 \quad (32)$$

が成立するとき、非線形大域結合システムは漸近安定と

なる。ここで、

$$\begin{aligned} \mathbf{M}(\mathbf{x}) &= \mathbf{P} \mathbf{G}(\mathbf{x}) (\boldsymbol{\Gamma}_\alpha - \boldsymbol{\Gamma}_\mu - \boldsymbol{\Gamma}_K) \\ &\quad + (\boldsymbol{\Gamma}_\alpha - \boldsymbol{\Gamma}_\mu - \boldsymbol{\Gamma}_K)^T \mathbf{G}(\mathbf{x}) \mathbf{P} \prec 0 \end{aligned} \quad (33)$$

であり、 \mathbf{P} は ii ブロックに正定対称行列 \mathbf{P}_i をもつブロック対角行列である。 $\boldsymbol{\Gamma}_K \in \mathbb{R}^{M \times M}$ は ii ブロックに $K_{ii} \boldsymbol{\Gamma}_{ii} \in \mathbb{R}^{n_i \times n_i}$ をもつブロック対角行列であり、 K_{ii} は $\boldsymbol{\Gamma}_K = \text{diag}\{(k_1 - 1)d_1, (k_2 - 1)d_2, (k_3 - 1)d_3, \dots, (k_m - 1)d_m, 0, 0, 0, \dots, 0\} \in \mathbb{R}^{N \times N}$ の ii 要素である。

(証明) 一般性を失うことなく $\mathbf{x}^* = \mathbf{0}$ を仮定する。

$$\begin{aligned} \dot{V}(\mathbf{x}) &= \nabla_{\mathbf{x}} V^T \mathbf{G}(\mathbf{x}) \{ \mathbf{f}(\mathbf{x}) + \boldsymbol{\Gamma}_\alpha \mathbf{x} - \boldsymbol{\Gamma}_K \mathbf{x} \} \\ &< \nabla_{\mathbf{x}} V^T \mathbf{G}(\mathbf{x}) (\boldsymbol{\Gamma}_\alpha - \boldsymbol{\Gamma}_\mu - \boldsymbol{\Gamma}_K) \mathbf{x} \\ &= 2\mathbf{x}^T \mathbf{P} \mathbf{G}(\mathbf{x}) (\boldsymbol{\Gamma}_\alpha - \boldsymbol{\Gamma}_\mu - \boldsymbol{\Gamma}_K) \mathbf{x} \\ &= \mathbf{x}^T \mathbf{M}(\mathbf{x}) \mathbf{x} \end{aligned} \quad (34)$$

よって、 $\mathbf{M}(\mathbf{x}) \prec 0$ で漸近安定となる。□

4.3 N 人進化ゲームの漸近安定化

N 人進化ゲームのダイナミクスにおいて、 $\mathbf{A}_{ij} = \mathbf{V}_{ij} - \frac{1}{N} \boldsymbol{\Theta}_{ij}$ とする。ただし、 $\boldsymbol{\Theta}_{ij} = [\theta_i \theta_i \theta_i \cdots \theta_i] \in \mathbb{R}^{n_i \times n_j}$ である。このとき、

$$\boldsymbol{\phi}_i = \sum_{j=1}^N \mathbf{A}_{ij} \mathbf{x}_j \quad (\forall i) \quad (35)$$

となる。ここで、

$$\mathbf{G}_{ii}(\mathbf{x}_i) = \mathbf{I}(\mathbf{x}_i) \{ \mathbf{I}_{n_i} - (\mathbf{1}_{n_i} \otimes \mathbf{x}_i^T) \} \quad (\forall i) \quad (36)$$

を考える。ただし、 $\mathbf{I}(\mathbf{x}_i) = \text{diag}\{x_i^1, x_i^2, x_i^3, \dots, x_i^{n_i}\} \in \mathbb{R}^{n_i \times n_i}$, $\mathbf{I}_{n_i} \in \mathbb{R}^{n_i \times n_i}$ は単位行列、 $\mathbf{1}_{n_i} = [1 \ 1 \ 1 \ \cdots \ 1]^T \in \mathbb{R}^{n_i \times 1}$ であり、 \otimes は Kronecker 積である。

いま、 N 人進化ゲームは、

$$\dot{\mathbf{x}}_i = \mathbf{G}(\mathbf{x}) \mathbf{A} \mathbf{x} \quad (37)$$

で表せることがわかる。 $\mathbf{A} \in \mathbb{R}^{M \times M}$ は ij ブロックに $\mathbf{A}_{ij} \in \mathbb{R}^{n_i \times n_j}$ をもつ。 $\mathbf{G}(\mathbf{x}) \in \mathbb{R}^{M \times M}$ は ii ブロックに $\mathbf{G}_{ii}(\mathbf{x}) \in \mathbb{R}^{n_i \times n_i}$ をもつブロック対角行列である。

ここで、与えられたペイオフのもとで均衡が漸近安定とならない場合に、ペイオフを変更することで漸近安定化する報酬の設計を考える。

【定理 2】 報酬の設計

ペイオフ \mathbf{A} による均衡 \mathbf{x}^* が不安定であるとき、 \mathbf{x}^* を漸近安定化する報酬の設計法は $\mathbf{A} + \bar{\mathbf{K}} \bar{\mathbf{P}}^{-1}$ により与えられる。

(証明) $\boldsymbol{\Gamma}_\alpha = \mathbf{A}$, $\mathbf{f}(\mathbf{x}) = \mathbf{0}$ として、

$$\begin{aligned} \mathbf{M}(\mathbf{x}^*) &= \mathbf{P} \mathbf{G}(\mathbf{x}^*) (\mathbf{A} - \boldsymbol{\Gamma}_K) \\ &\quad + (\mathbf{A} - \boldsymbol{\Gamma}_K)^T \mathbf{G}(\mathbf{x}^*) \mathbf{P} \prec 0 \end{aligned} \quad (38)$$

より、正定対称行列 \mathbf{Q} を用いると、

$$\begin{aligned} & \mathbf{G}(\mathbf{x}^*)\mathbf{A}\bar{\mathbf{P}} + \mathbf{G}(\mathbf{x}^*)\bar{\mathbf{K}} \\ & + (\mathbf{G}(\mathbf{x}^*)\mathbf{A}\bar{\mathbf{P}} + \mathbf{G}(\mathbf{x}^*)\bar{\mathbf{K}})^T + \mathbf{Q} \preceq 0 \end{aligned} \quad (39)$$

となる。 $\bar{\mathbf{P}} = \mathbf{P}^{-1}$, $\bar{\mathbf{K}} = \Gamma_K \mathbf{P}^{-1}$ と \mathbf{Q} に関する線形行列不等式より、 $\bar{\mathbf{P}}$ と $\bar{\mathbf{K}}$ が得られれば、 $\mathbf{A} - \Gamma_K = \mathbf{A} + \bar{\mathbf{K}}\bar{\mathbf{P}}^{-1}$ で設計できる。□

上記の遺伝子、神経回路網や進化ゲームは、プレイヤー i が戦略 π を選び行動した結果、ペイオフにもとづく効用関数の変化を入力として動作していると考えると、プレイヤー j と戦略 π' に対するペイオフ $v_{ij}^{\pi, \pi'} (\forall j, \forall \pi')$ すべて既知とする報酬駆動型システムであるといえる。

5. 報酬による最適化

5.1 遺伝的アルゴリズム

複製方程式で表現できる遺伝子のモデルでは、淘汰と突然変異が考慮されていた。ここでは、さらに交叉と確率を考慮した遺伝子のモデルである以下のような簡易な遺伝的アルゴリズム (Simple Genetic Algorithm) について考える [9]。

- [1] l ビットの 0 または 1 で遺伝子を表現する。 N 個の個体 $\{\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \dots, \mathbf{x}_N\}$ を作成。 $N \ll 2^l - 1 = n$ とする。
- [2] 適合度関数 $f(\cdot)$ 、突然変異率 ϵ 、 α 、交叉率 $\chi_{\mathbf{k}}$ の値と逆温度 β の初期値を与える。
- [3] 突然変異する部分 \mathbf{h} と交叉する部分 \mathbf{k} を与える。
- [4] 選択、突然変異、交叉を繰り返し世代交代する。

淘汰は、現世代において存在する遺伝子群 $\mathbf{p} = \{\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \dots, \mathbf{x}_N\}$ のうち、遺伝子が選択される確率

$$S_{\beta}(\mathbf{x}'|\mathbf{p}) = \frac{f(\mathbf{x}')^{\beta}}{\sum_{\mathbf{x} \in \mathbf{p}} f(\mathbf{x})^{\beta}} \quad (\mathbf{x}' \in \mathbf{p}) \quad (40)$$

で表現される。ここで、 $f(\mathbf{x})$ は遺伝子 \mathbf{x} に対する適応度関数であり、 $\beta > 0$ は Boltzmann の逆温度である。

突然変異は、確率 $\mu_{\beta} = \epsilon \exp(-\alpha\beta)$ によりビットが反転する確率

$$m_{\beta}(\mathbf{x}, \mathbf{h}) = \mu_{\beta}^{\mathbf{1}^T \mathbf{h}} (1 - \mu_{\beta})^{l - \mathbf{1}^T \mathbf{h}} \quad (41)$$

で表現される。ここで、 $\mathbf{1} \in \mathbb{R}^{l \times 1}$ は要素すべてが 1 のベクトルであり、 $\mathbf{h} \in \mathbb{R}^{l \times 1}$ は突然変異で反転する部分のビットは 1、それ以外は 0 で与えられる。 $\mathbf{x}' = \mathbf{x} \oplus \mathbf{h}$ が突然変異後の状態を表す。ただし、 \oplus は排他的論理和である。

交叉は、遺伝子 \mathbf{x}_i と \mathbf{x}_j が交叉率 $\chi_{\mathbf{k}}$ で交叉することで、遺伝子 \mathbf{x}' となる確率

$$\begin{aligned} C_{\mathbf{x}'}(\mathbf{x}_i, \mathbf{x}_j) &= \sum_{\mathbf{k}} \frac{\chi_{\mathbf{k}} + \chi_{\bar{\mathbf{k}}}}{2} \\ &\times \delta[\mathbf{x}_i \odot \mathbf{k} \oplus \bar{\mathbf{k}} \odot \mathbf{x}_j = \mathbf{x}'] \end{aligned} \quad (42)$$

で表現される。ここで、 $\mathbf{k} \in \mathbb{R}^{l \times 1}$ は交叉する部分のビットは 1、それ以外は 0 で与えられる。 $\bar{\mathbf{k}}$ は \mathbf{k} の否定である。

り、 $\delta[\cdot]$ は真なら 1、それ以外は 0 となる条件演算子である。ただし、 \odot は論理積である。

遺伝子 \mathbf{x} が淘汰、突然変異、交叉を経て遺伝子 \mathbf{x}' となる遷移確率 $\mathcal{M}(\mathbf{x}'|\mathbf{x})$ は

$$\begin{aligned} & \mathcal{M}(\mathbf{x}'|\mathbf{x}) \\ &= \sum_{\mathbf{y}, \mathbf{h}_1, \mathbf{h}_2} C_{\mathbf{x}'}(m_{\beta}(S_{\beta}(\mathbf{x}|\mathbf{p}), \mathbf{h}_1), m_{\beta}(S_{\beta}(\mathbf{y}|\mathbf{p}), \mathbf{h}_2)) \\ &= S_{\beta}(\mathbf{x}|\mathbf{p}) \sum_{\mathbf{y}} S_{\beta}(\mathbf{y}|\mathbf{p}) \\ &\quad \times \sum_{\mathbf{h}_1, \mathbf{h}_2} m_{\beta}(\mathbf{x}, \mathbf{h}_1) m_{\beta}(\mathbf{y}, \mathbf{h}_2) \sum_{\mathbf{k}} \frac{\chi_{\mathbf{k}} + \chi_{\bar{\mathbf{k}}}}{2} \\ &\quad \times \delta[(\mathbf{x} \oplus \mathbf{h}_1) \odot \mathbf{k} \oplus \bar{\mathbf{k}} \odot (\mathbf{y} \oplus \mathbf{h}_2) = \mathbf{x}'] \end{aligned} \quad (43)$$

となる。

5.2 報酬駆動型システムとしての遺伝的アルゴリズム

遷移確率 $\mathcal{M}(\mathbf{x}'|\mathbf{x})$ が遺伝子 \mathbf{x} から遺伝子 \mathbf{x}' へ切り替るために必要な制御入力 $\phi_{\mathbf{x}', \mathbf{x}}$

$$\mathcal{M}(\mathbf{x}'|\mathbf{x}) = p(\mathbf{x}'|\mathbf{x}) \exp(\phi_{\mathbf{x}', \mathbf{x}}) \quad (44)$$

によりモデル化され、遺伝的アルゴリズムが累積無限期間平均コスト関数

$$\mathcal{R}^{\pi}(\mathbf{x}) = \lim_{t_f \rightarrow \infty} \mathbb{E} \left[\sum_{\tau=0}^{t_f-1} l(\mathbf{x}_{(\tau)}, \pi(\mathbf{x}_{(\tau)})) \right] \quad (45)$$

の最小化 $\mathcal{R}(\mathbf{x}) = \min_{\pi} \mathcal{R}^{\pi}(\mathbf{x})$ を評価していると考える。 $p(\mathbf{x}'|\mathbf{x})$ は受動的遷移であり、 $p(\mathbf{x}'|\mathbf{x}) = 0$ なら $\mathcal{M}(\mathbf{x}'|\mathbf{x}) = 0$ である。また、 $\sum_{\mathbf{x}'} \mathcal{M}(\mathbf{x}'|\mathbf{x}) = \sum_{\mathbf{x}'} p(\mathbf{x}'|\mathbf{x}) = 1$ 、 $\pi(\mathbf{x}_{(\tau)}) = \phi_{\mathbf{x}_{(\tau)}'; \mathbf{x}_{(\tau)}}$ である。

コスト関数 $l(\mathbf{x}, \mathbf{u})$ は報酬 $r(\mathbf{x}'; \mathbf{x})$ と制御入力 $\phi_{\mathbf{x}', \mathbf{x}}$ を用いて

$$l(\mathbf{x}, \phi) = \mathbb{E}_{\mathbf{x}' \sim \mathcal{M}(\cdot | \mathbf{x})} [r(\mathbf{x}'|\mathbf{x}) + \phi_{\mathbf{x}', \mathbf{x}}] \quad (46)$$

で定義されるものとする。

このとき、Bellman 方程式は

$$\begin{aligned} \mathcal{R}(\mathbf{x}) &= \min_{\phi \in \phi_{\mathbf{x}', \mathbf{x}}} \{l(\mathbf{x}, \mathbf{u}) \\ &\quad + \mathbb{E}_{\mathbf{x}' \sim \mathcal{M}(\cdot | \mathbf{x})} [\mathcal{R}(\mathbf{x}')] \} - c \end{aligned} \quad (47)$$

となる。ここで、 c は平均コストである。

最適制御入力は、 $\psi(\mathbf{x}'; \mathbf{x}) = -r(\mathbf{x}'; \mathbf{x}) - \mathcal{R}(\mathbf{x}')$ として

$$\phi_{\mathbf{x}', \mathbf{x}}^* = \psi(\mathbf{x}'; \mathbf{x}) - \log \mathbb{E}_{\mathbf{x}' \sim p(\cdot | \mathbf{x})} [\exp(\psi(\mathbf{x}'; \mathbf{x}))] \quad (48)$$

となる [10]。

ここで、逐次得られる、限られたペイオフである適応度のもとで、最も適応度が高い遺伝子の表現を探索する報酬による最適化を考える。

【定理 3】 報酬による最適化

簡易な遺伝的アルゴリズムでは、報酬 $r(\mathbf{x}'; \mathbf{x})$ を遺伝

子 \mathbf{x} から遺伝子 \mathbf{x}' へ戦略を切り替えたときの適応度の変化

$$r(\mathbf{x}'; \mathbf{x}) = \log p(\mathbf{x}' | \mathbf{x}) \left(\frac{f(\mathbf{x})}{f(\mathbf{x}')} \right)^\beta \quad (49)$$

で与えて駆動することで、累積無限期間平均コスト関数

$$\mathcal{R}^\pi(\mathbf{x}) = \lim_{t_f \rightarrow \infty} \mathbb{E} \left[\sum_{\tau=0}^{t_f-1} \mathbb{E}_{\mathbf{x}' \sim \mathcal{M}(\cdot | \mathbf{x})} \left[\log \left(\frac{f(\mathbf{x}_{(\tau)})}{f(\mathbf{x}'_{(\tau)})} \right) \right] \right] \quad (50)$$

の最小化を行っている。

(証明) コスト関数 $l(\mathbf{x}, \phi)$ は

$$\begin{aligned} l(\mathbf{x}, \phi) &= \mathbb{E}_{\mathbf{x}' \sim \mathcal{M}(\cdot | \mathbf{x})} [v(\mathbf{x}' | \mathbf{x}) + \phi \mathbf{x}' \cdot \mathbf{x}] \\ &= \beta \mathbb{E}_{\mathbf{x}' \sim \mathcal{M}(\cdot | \mathbf{x})} \left[\log \left(\frac{f(\mathbf{x})}{f(\mathbf{x}') \exp(\phi \mathbf{x}' \cdot \mathbf{x})} \right) \right] - S_{\mathcal{M}} \end{aligned} \quad (51)$$

とできる。ここで、 $S_{\mathcal{M}} = \mathbb{E}_{\mathbf{x}' \sim \mathcal{M}(\cdot | \mathbf{x})} [-\log \mathcal{M}(\mathbf{x}' | \mathbf{x})]$ はエントロピーである。□

ところで、最適制御入力 $\phi_{\mathbf{x}', \mathbf{x}}^*$ は、以下の Kullback-Leibler divergence

$$\text{KL}(\mathcal{M}(\mathbf{x}' | \mathbf{x}) \| p(\mathbf{x}' | \mathbf{x}) \exp(\phi_{\mathbf{x}', \mathbf{x}})) = \mathbb{E}_{\mathbf{x}' \sim \mathcal{M}(\cdot | \mathbf{x})} \left[\frac{\mathcal{M}(\mathbf{x}' | \mathbf{x})}{p(\mathbf{x}' | \mathbf{x}) \exp(\phi_{\mathbf{x}', \mathbf{x}})} \right] \quad (52)$$

を最小化している。遷移確率のモデル化において、Kolmogorov-Gabor 多項式で表現されるポテンシャルの存在を仮定する場合には、ボルツマンマシンと類似の手法が適用可能である [11]。

遺伝的アルゴリズムは、戦略を選び行動した結果、逐次、ペイオフである適応度が判明し、報酬にもとづき入力が得られる報酬駆動型システムであるといえる。

6. おわりに

本稿では、意思決定の主体が戦略を選び行動した結果、ペイオフにもとづく効用関数の変化を報酬と考え、それを入力として動作する報酬駆動型システムを考えた。まず、効用関数の最大化が適者生存則に関連していることを紹介した。つぎに、与えられたペイオフのもとで均衡が漸近安定とならない場合に、ペイオフを変更することで漸近安定化する報酬の設計法を概説した。そして、逐次得られる限られたペイオフのもとで、最もペイオフが高い戦略を探査する報酬による最適化について概観した。

(2014年6月2日受付)

参考文献

- [1] J. Hofbauer and K. Sigmund: *The Theory of Evolution and Dynamical Systems*, Cambridge University Press (1988)

- [2] T. Yamada and H. Fujisaka: Stability theory of synchronized motion in coupled-oscillator systems. II –The mapping approach–; *Progress of Theoretical Physics*, Vol. 70, No. 5, pp. 1240–1248 (1983)
- [3] 金子, 津田: 複雑系のカオス的シナリオ, 朝倉書店, p. 127 (1996)
- [4] 堀江, 相吉, 宮野: ゲーム理論的均衡解探索のための疑似勾配系モデルのニューラルネットワーク実現とその運動; システム/制御/情報, Vol. 12, No. 11, pp. 680–690 (1999)
- [5] 奥原, 尾崎: 適者生存型学習則を適用した競合動径基底関数ネットワーク; 電子情報通信学会論文誌, Vol. J80-DII, No. 12, pp. 3191–3199 (1997)
- [6] K. Okuhara, K. Sasaki and S. Osaki: Reproductive and competitive radial basis function networks adaptable to dynamical environments; *Systems and Computers in Japan*, Vol. 31, No. 13, pp. 65–73 (2000)
- [7] J. Xiang and G. Chen: On the V-stability of complex dynamical networks; *Automatica*, Vol. 43, pp. 1049–1057 (2007)
- [8] X. Li, X. Wang and G. Chen: Pinning a complex dynamical network to its equilibrium; *IEEE Transactions on Circuits and Systems*, Vol. 51, No. 10, pp. 2074–2087 (2004)
- [9] M. D. Vose: *The Simple Genetic Algorithm: Foundations and Theory*, MIT Press, (1999)
- [10] E. Todorov: Efficient computation of optimal actions; *PNAS*, Vol. 106, No. 28, pp. 11478–11483 (2009)
- [11] K. Okuhara, S. Osaki and M. Kijima: Learning to design synergetic computers with an expanded symmetric diffusion network; *Neural Computation*, Vol. 11, pp. 1475–1491 (1999)

著者略歴

奥原 浩之 (正会員)



1996年3月広島大学大学院工学研究科システム工学専攻博士後期課程修了。同年4月同大学工学部助手、1998年10月広島県立大学経営学部講師、2000年同助教授、2006年4月大阪大学大学院情報科学研究科准教授となり現在に至る。2006年4月大阪大学金融保険教育センター兼任、2007年4月台湾国立成功大学客員准教授、2010年4月大学共同利用機関法人情報・システム研究機構統計数理研究所客員准教授。ソフトコンピューティング技術の開発、数理モデリングと最適化・制御の研究に従事。2006年スケジューリング学会技術賞受賞。IEEE、電子情報通信学会、日本オペレーションズ・リサーチ学会などの会員。