



研究ノート

報酬駆動型システム： 自律分散メカニズム・デザイン実現のための研究

奥 原 浩 之*

Reward Derived System:
Approach for Autonomous Decentralized Mechanism Design

Key Words : Reward Structures, Autonomous Decentralized, Design, Optimization

1 はじめに

20世紀後半頃からの地球の温暖化が意識され、温室効果ガスの排出抑制が国際的に議論されて久しい。地球温暖化が気候に影響するのなら、生態系などの自然環境が変化し、農業や漁業を通じて社会環境へも影響が及ぶ可能性が考えられている。実効性がある温暖化抑制策あるいは緩和策の実施の遅れは、さらなる負担やコストの増加を招くと指摘されている。にもかかわらず、利害の関係がまちまちであるため制度の設計は困難しているように見受けられる。

これに類似する問題解決の重要性は、なにも世界規模の事例を取り上げるまでもなく、国内、地域社会といった規模を問わず、我々の身近な問題で普段感じられるものであり、制度の設計の困難さを示していると言える。

常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とは考えられない。そのような場合には、環境を人為的に設計することで、個体の最適な応答により、集団として望ましい効率的な選択を自律分散的に実現する仕組みの解明や設計法が求められる。

本研究ノートでは、特定の戦略の選択を促進・抑制する自律分散メカニズム・デザインを実現するための研究[1]の一部を簡単に紹介する。文献[1]では遺伝的アルゴリズム、神経回路網モデル、進化ゲ

ームに関する力学は、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報酬系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事が示されているが、ここでは、遺伝的アルゴリズムについて述べる。

2 遺伝的アルゴリズム

淘汰と突然変異、さらに交叉と確率を考慮した遺伝子のモデルである簡易な遺伝的アルゴリズム(Simple Genetic Algorithm:SGA)[2]は次のようなものである。個体の最適な応答により、集団として望ましい効率的な選択を自律分散的に実現する仕組みとなっている。

- [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 $\{x_1, x_2, x_3, \dots, x_N\}$ 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、交叉率の値と逆温度 β の初期値を与える。
- [3] 突然変異する部分と交叉する部分を与える。
- [4] 選択、突然変異、交叉を繰り返し世代交代する。

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

淘汰は、現世代において存在する遺伝子群 $p = \{x_1, x_2, x_3, \dots, x_N\}$ のうち、遺伝子が選択される確率で表現される。突然変異は、ビットが反転する確率で表現される。

交叉は、遺伝子 x_i と x_j が交叉率で交叉することで、遺伝子 x' となる確率で表現される。

3 報酬系としての遺伝的アルゴリズム

遺伝子 x が淘汰、突然変異、交叉を経て遺伝子



* Koji OKUHARA

1968年9月生
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（1996年）
現在、大阪大学 大学院情報科学研究科
准教授 博士（工学）
ソフトコンピューティング、数理モデリングと最適化・制御
TEL：06-6879-7877
FAX：06-6879-7877
E-mail：okuhara@ist.osaka-u.ac.jp

x' となる遷移確率 $M(x'|x)$ が

$$M(x'|x) = p(x'|x) \exp(\phi_{x';x})$$

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

$$R^\pi(x) = \lim_{t_f \rightarrow \infty} E \left[\sum_{\tau=0}^{t_f-1} l(x_{(\tau)}, \pi(x_{(\tau)})) \right]$$

の最小化 $R(x) = \min_\pi R^\pi(x)$ を評価していると考える。ここで、 $\phi_{x';x}$ は遺伝子 x から遺伝子 x' へ切替えるために必要な制御入力であり、 $p(x'|x)$ は受動的遷移であり、 $p(x'|x) = 0$ なら $M(x'|x) = 0$ である。

また、 $\sum_{x'} M(x'|x) = \sum_{x'} p(x'|x) = 1$ 、 $\pi(x_{(\tau)}) = \phi_{x'_{(\tau)};x_{(\tau)}}$ である。

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

$$l(x, \phi) = E_{x' \sim M(\cdot|x)} [r(x'; x) + \phi_{x';x}]$$

で定義されるものとする。

このとき、Bellman 方程式は

$$\begin{aligned} R(x) &= \min_{\phi \in \phi_{x';x}} \{ l(x, u) \\ &\quad + E_{x' \sim M(\cdot|x)} [R(x')] \} - c \end{aligned}$$

となる。ここで、 c は平均コストである。

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

$$\phi_{x';x}^* = \psi(x'; x) - \log E_{x' \sim p(\cdot|x)} [\exp(\psi(x'; x))]$$

となる [3]。

このことは、簡易な遺伝的アルゴリズムでは、報酬 $r(x'; x)$ を遺伝子 x から遺伝子 x' へ戦略を切替えたときの適応度の変化

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

で与えて駆動することで、累積無限期間平均コスト関数

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

の最小化を行っていることを示している。

報酬駆動型システムとしての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プレイヤー間で共有され、報酬にもとづき入力が得られる報酬系であることがわかる。つまり、個体の最適な応答により、集団として望ましい効率的な選択を自律分散的に実現する仕組みを、報酬系による最適化を行うことで実現していると言える。

4 おわりに

本研究ノートでは、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による報酬系を考えた。その枠組みの中で、遺伝的アルゴリズムが報酬系による最適化を行っており、自律分散メカニズム・デザインを実現していることを明らかにした。

参考文献

- [1] 奥原浩之, 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, システム / 制御 / 情報, Vol. 58, No. 11, pp. 462-467 (2014)
- [2] M. D. Vose, *The Simple Genetic Algorithm: Foundations and Theory*, MIT Press (1999)
- [3] E. Todorov, Efficient Computation of Optimal Actions; PNAS, Vol. 106, No. 28, pp.11478-11483 (2009)