

マルコフ決定過程に基づくマルチエージェントシステムの漸近的性質

An Asymptotic Property of Multiagent Systems Based on a Markov Decision Process

岩田 一貴†
Iwata Kazunori

池田 和司‡
Ikeda Kazushi

酒井 英昭‡
Sakai Hideaki

1 はじめに

複数のエージェントが、互いに協調することにより、ある与えられたタスクに関して最適な政策を学習するようなマルチエージェントシステム (MAS) を考える。MAS が従う確率過程 (枠組み) としては、Temporal Difference (TD) 学習などの代表的な強化学習アルゴリズムが適用し易いなどの理由から、マルコフ決定過程 (MDP) がよく知られている (例えば, [1] を参照)。本論では, MAS が定常エルゴード MDP に従う場合に成立する漸近的性質を示し, 最適政策の学習により MAS の収益が最大化される過程を確率的な観点から明らかにする。

2 準備

初めに記号の準備をする。今, 同じ環境に存在する任意の M 人のエージェントに 1 から M の番号をつけ, 番号の集合を $\mathcal{M} \triangleq \{1, \dots, M\}$ とする。MAS の状態の集合を $\mathcal{S} \triangleq \{s_1, \dots, s_I\}$, エージェントがとり得る行動の集合を $\mathcal{A} \triangleq \{a_1, \dots, a_J\}$, 実数の集合 \mathbb{R} を適当に離散化した集合を $\mathbb{R}_0 \triangleq \{r_1, \dots, r_K\} \subset \mathbb{R}$ とする。よって, $|\mathcal{S}| = I$, $|\mathcal{A}| = J$, $|\mathbb{R}_0| = K$ である。記号の簡単のために, MAS 全体 (M 人のエージェント) のある状態を $\mathbf{s}_{i_1 \dots i_M} = (s_{i_1}, \dots, s_{i_M}) \in \mathcal{S}^M$ と書く。ここで, m 番目の要素 s_{i_m} ($i_m \in \{1, \dots, I\}$) はエージェント $m \in \mathcal{M}$ の状態を示す。同様に, エージェント全体のある行動を $\mathbf{a}_{j_1 \dots j_M} = (a_{j_1}, \dots, a_{j_M}) \in \mathcal{A}^M$, ある報酬を $\mathbf{r}_{k_1 \dots k_M} = (r_{k_1}, \dots, r_{k_M}) \in \mathbb{R}_0^M$ と書く。変数 $t \in \mathbb{N}$ は時間ステップを表し, $\mathbf{s}_m(t)$, $\mathbf{a}_m(t)$, $r_m(t)$ はそれぞれ時間ステップ t におけるエージェント $m \in \mathcal{M}$ の状態, 実行した行動, 観測した報酬を示す確率変数とする。同様に, 時間ステップ t における MAS 全体の状態, 行動, 報酬の確率変数をそれぞれ $\mathbf{s}(t) \triangleq \{s_m(t) | m \in \mathcal{M}\}$, $\mathbf{a}(t) \triangleq \{a_m(t) | m \in \mathcal{M}\}$, $\mathbf{r}(t) \triangleq \{r_m(t) | m \in \mathcal{M}\}$ と表す。

2.1 MAS の定常エルゴード MDP

定常エルゴード MDP において, ある初期状態分布 $\{\Pr(\mathbf{s}(1) = \mathbf{s}_{i_1 \dots i_M}) | \mathbf{s}_{i_1 \dots i_M} \in \mathcal{S}^M\}$ が与えられたとき, MAS の確率変数 $\mathbf{s}(t)$, $\mathbf{a}(t)$, $\mathbf{r}(t)$ に関する確率分布は, エージェントの政策行列

$$\Gamma^\pi \triangleq (\mathbf{P}_{(11 \dots 1)}, \dots, \mathbf{P}_{(II \dots I)})^\top, \quad (1)$$

および環境の状態遷移行列 Γ^T

$$\Gamma^T \triangleq (\mathbf{P}_{(11 \dots 1, 11 \dots 1)}, \dots, \mathbf{P}_{(II \dots I, JJ \dots J)})^\top, \quad (2)$$

で定義される。ただし, Γ^π は $I^M \times J^M$ 行列で, 各要素を

$$\mathbf{P}_{(i_1 \dots i_M)} \triangleq (p_{i_1 \dots i_M, 11 \dots 1}, \dots, p_{i_1 \dots i_M, JJ \dots J})^\top, \quad (3)$$

と書く。同様に, Γ^T は $I^M J^M \times I^M K^M$ 行列で, 各要素を

$$\mathbf{P}_{(i_1 \dots i_M, j_1 \dots j_M)} \triangleq (p_{i_1 \dots i_M, j_1 \dots j_M, 11 \dots 1, 11 \dots 1}, \dots, p_{i_1 \dots i_M, j_1 \dots j_M, II \dots I, KK \dots K})^\top, \quad (4)$$

と記述する。さらに, (3) の各要素の確率は,

$$p_{i_1 \dots i_M, j_1 \dots j_M} \triangleq \Pr(\mathbf{a}(t) = \mathbf{a}_{j_1 \dots j_M} | \mathbf{s}(t) = \mathbf{s}_{i_1 \dots i_M}), \quad (5)$$

を示し, (4) の各要素の確率は,

$$p_{i_1 \dots i_M, j_1 \dots j_M, i'_1 \dots i'_M, k_1 \dots k_M} \triangleq \Pr(\mathbf{s}(t+1) = \mathbf{s}'_{i'_1 \dots i'_M}, \mathbf{r}(t+1) = \mathbf{r}_{k_1 \dots k_M} | \mathbf{s}(t) = \mathbf{s}_{i_1 \dots i_M}, \mathbf{a}(t) = \mathbf{a}_{j_1 \dots j_M}), \quad (6)$$

を示す。本来は各エージェントが学習する場合, 政策行列 Γ^π は時変となる。しかし, 簡単のために, 本論では政策行列 Γ^π が十分にゆっくりと学習により更新されることを仮定し, MAS が定常エルゴード MDP に従うと見なす。このような仮定に基づく解析は, TD 学習などの確率近似法に基づく学習に対してよく行われる (例えば, [2] を参照)。また, 環境の状態遷移行列 Γ^T の各要素は定数で, 各エージェントは Γ^T がわからないものとする。

2.2 定常エルゴード MDP の経験分布

任意の n に対する定常エルゴード MDP の経験系列 $\mathbf{x} \in (\mathcal{S} \times \mathcal{A} \times \mathbb{R}_0)^{nM}$ を確率変数列

$$\mathbf{s}(1), \mathbf{a}(1), \mathbf{s}(2), \mathbf{r}(2), \mathbf{a}(2), \dots, \mathbf{s}(n), \mathbf{r}(n), \mathbf{a}(n), \mathbf{r}(n+1),$$

の観測値とする。ある経験系列 $\mathbf{x} \in (\mathcal{S} \times \mathcal{A} \times \mathbb{R}_0)^{nM}$ が与えられたとき, 任意の状態 $\mathbf{s}_{i_1 \dots i_M} \in \mathcal{S}^M$, 状態-行動 $(\mathbf{s}_{i_1 \dots i_M}, \mathbf{a}_{j_1 \dots j_M}) \in \mathcal{S}^M \times \mathcal{A}^M$, 状態-行動-次状態-報酬 $(\mathbf{s}_{i_1 \dots i_M}, \mathbf{a}_{j_1 \dots j_M}, \mathbf{s}'_{i'_1 \dots i'_M}, \mathbf{r}_{k_1 \dots k_M}) \in \mathcal{S}^M \times \mathcal{A}^M \times \mathcal{S}^M \times \mathbb{R}_0^M$ の出現回数をそれぞれ

$$n_{i_1 \dots i_M} \triangleq |\{t | \mathbf{s}(t) = \mathbf{s}_{i_1 \dots i_M}, t = 1, \dots, n\}|, \quad (7)$$

$$n_{i_1 \dots i_M, j_1 \dots j_M} \triangleq |\{t | (\mathbf{s}(t), \mathbf{a}(t)) = (\mathbf{s}_{i_1 \dots i_M}, \mathbf{a}_{j_1 \dots j_M}), t = 1, \dots, n\}|, \quad (8)$$

† 広島市立大学情報科学部 〒731-3194 広島市安佐南区大塚東 3-4-1 Email: kiwata@im.hiroshima-cu.ac.jp
‡ 京都大学大学院情報学研究所 〒606-8501 京都市左京区吉田本町
Email: {kazushi, hsakai}@i.kyoto-u.ac.jp

$$n_{i_1 \dots i_M, j_1 \dots j_M, i'_1 \dots i'_M, k_1 \dots k_M} \triangleq |\{t \mid (s(t), \mathbf{a}(t), \mathbf{s}(t+1), \mathbf{r}(t+1)) = (s_{i_1 \dots i_M}, \mathbf{a}_{j_1 \dots j_M}, \mathbf{s}'_{i'_1 \dots i'_M}, \mathbf{r}_{k_1 \dots k_M}), t = 1, \dots, n)\}|, \quad (9)$$

と表す．ただし，(9) では $s(n+1) = s(1)$ と約束する．このとき，MAS の状態 $s_{i_1 \dots i_M}$ に関する経験分布 $f_{i_1 \dots i_M} \in [0, 1]$ (一般にタイプと呼ぶ [3]) と状態-行動 $(s_{i_1 \dots i_M}, \mathbf{a}_{j_1 \dots j_M})$ に関する同時タイプ $f_{i_1 \dots i_M, j_1 \dots j_M} \in [0, 1]$ は

$$n_{i_1 \dots i_M} = n f_{i_1 \dots i_M}, \quad (10)$$

$$n_{i_1 \dots i_M, j_1 \dots j_M} = n f_{i_1 \dots i_M, j_1 \dots j_M}, \quad (11)$$

により定義される．同様に，状態-行動 $(s_{i_1 \dots i_M}, \mathbf{a}_{j_1 \dots j_M})$ に関する条件付タイプ $g_{i_1 \dots i_M, j_1 \dots j_M} \in [0, 1]$ は，

$$n_{i_1 \dots i_M, j_1 \dots j_M} = n_{i_1 \dots i_M} g_{i_1 \dots i_M, j_1 \dots j_M}, \quad (12)$$

状態-行動-次状態-報酬 $(s_{i_1 \dots i_M}, \mathbf{a}_{j_1 \dots j_M}, \mathbf{s}'_{i'_1 \dots i'_M}, \mathbf{r}_{k_1 \dots k_M})$ に関するマルコフタイプ $g_{i_1 \dots i_M, j_1 \dots j_M, i'_1 \dots i'_M, k_1 \dots k_M} \in [0, 1]$ は，

$$n_{i_1 \dots i_M, j_1 \dots j_M, i'_1 \dots i'_M, k_1 \dots k_M} = n_{i_1 \dots i_M, j_1 \dots j_M} g_{i_1 \dots i_M, j_1 \dots j_M, i'_1 \dots i'_M, k_1 \dots k_M}, \quad (13)$$

により定義される．簡単のために，(10)–(13) で定義されるタイプを要素に持つ行列 (もしくはベクトル) をそれぞれ

$$\mathbf{F}_{S^M} \triangleq \{f_{i_1 \dots i_M}\}, \quad (14)$$

$$\mathbf{F}_{S^M A^M} \triangleq \{f_{i_1 \dots i_M, j_1 \dots j_M}\}, \quad (15)$$

$$\Phi^\pi \triangleq \{g_{i_1 \dots i_M, j_1 \dots j_M}\}, \quad (16)$$

$$\Phi^\pi \triangleq \{g_{i_1 \dots i_M, j_1 \dots j_M, i'_1 \dots i'_M, k_1 \dots k_M}\}, \quad (17)$$

と定義する．さらに，状態 $s_{i_1 \dots i_M}$ および状態-行動 $(s_{i_1 \dots i_M}, \mathbf{a}_{j_1 \dots j_M})$ に関する定常確率 (定常エルゴードならば存在は明らか) を

$$\mathbf{V} \triangleq \lim_{n \rightarrow \infty} \mathbf{F}_{S^M}, \quad \mathbf{W} \triangleq \lim_{n \rightarrow \infty} \mathbf{F}_{S^M A^M}, \quad (18)$$

と定義する． $I^M \times J^M$ 行列 Φ^π および $I^M J^M \times I^M K^M$ 行列 Φ^π は，それぞれ政策確率行列 Γ^π および状態遷移行列 Γ^π に関するタイプを要素にもつ行列であることに注意されたい．

3 MAS の定常エルゴード MDP における漸近的性質

強化学習などの報酬を基にした教師無し学習の枠組みにおいて，MAS の最適政策は収益 (報酬の重み和) を最大にするような政策として定義される [1, 4]．よって，確率的な観点から見た場合，MAS の収益最大化とは経験系列に関する行列 Φ^π が最適な政策行列 (多くの場合，決定論的) に収束することである．

定義 1 (収益最大化) 任意の最適な政策行列 Γ^{π^\dagger} に対して，MAS の経験系列 $\mathbf{x} \in (\mathcal{S} \times \mathcal{A} \times \mathbb{R}_0)^{nM}$ に関する行列 Φ^π が

$$\lim_{n \rightarrow \infty} \Pr(\Phi^\pi = \Gamma^{\pi^\dagger}) = 1, \quad (19)$$

を満たすとき，確率収束という意味での収益最大化と呼ぶ．

この定義は [4, Def. 4] に相当する．Sanov の定理を定常エルゴード MDP に当てはめると， n に関する指数的な収束の速さ

$$\lim_{n \rightarrow \infty} \frac{1}{n} \log \Pr(\Phi^\pi = \Gamma^{\pi^\dagger}) = -D(\Gamma^{\pi^\dagger} \parallel \Gamma^\pi | \mathbf{V}), \quad (20)$$

が得られ，明らかに $n \rightarrow \infty$ に対して $\Gamma^\pi \rightarrow \Gamma^{\pi^\dagger}$ となれば収益最大化は達成されることがわかる．ただし，関数 D は

$$D(\Gamma^{\pi^\dagger} \parallel \Gamma^\pi | \mathbf{V}) = \sum_{i_1 \dots i_M} v_{i_1 \dots i_M} \times \sum_{i_1 \dots i_M, j_1 \dots j_M} p_{i_1 \dots i_M, j_1 \dots j_M}^\dagger \log \frac{p_{i_1 \dots i_M, j_1 \dots j_M}^\dagger}{p_{i_1 \dots i_M, j_1 \dots j_M}}, \quad (21)$$

を示し， $v_{i_1 \dots i_M}$ は \mathbf{V} の要素， $p_{i_1 \dots i_M, j_1 \dots j_M}^\dagger$ は Γ^{π^\dagger} の要素を表す．今，別の形で (20) を詳しく考察する．漸近等分割性より，

$$\lim_{n \rightarrow \infty} \Pr(\Phi^\pi = \Gamma^\pi, \Phi^\pi = \Gamma^\pi) = 1, \quad (22)$$

$$\lim_{n \rightarrow \infty} \frac{1}{n} \log P(\mathbf{x}) = -H(\Gamma^\pi | \mathbf{V}) - H(\Gamma^\pi | \mathbf{W}), \quad (23)$$

$$\lim_{n \rightarrow \infty} \frac{1}{n} \log |\{\mathbf{x} \in (\mathcal{S} \times \mathcal{A} \times \mathbb{R}_0)^{nM} \mid \Phi^\pi = \Gamma^\pi, \Phi^\pi = \Gamma^\pi\}| = H(\Gamma^\pi | \mathbf{V}) + H(\Gamma^\pi | \mathbf{W}), \quad (24)$$

が成り立つ．ただし，関数 H はエントロピーを表し，例えば (23) における $H(\Gamma^\pi | \mathbf{V})$ は， \mathbf{V} の要素を $v_{i_1 \dots i_M}$ とすると

$$H(\Gamma^\pi | \mathbf{V}) = \sum_{i_1 \dots i_M} v_{i_1 \dots i_M} \times \sum_{i_1 \dots i_M, j_1 \dots j_M} p_{i_1 \dots i_M, j_1 \dots j_M} \log \frac{1}{p_{i_1 \dots i_M, j_1 \dots j_M}}, \quad (25)$$

を示す．タイプ理論 [3] により，(22)–(24) は詳細な形で [4, Theo. 1–3] に示されているので証明は割愛する．任意の学習により $\Gamma^\pi \approx \Gamma^{\pi^\dagger}$ となるならば，(22)–(24) より有限だが十分大きな時間ステップ n における収益最大化の確率は，典型集合に含まれる $\mathcal{X}_n^\dagger = \{\mathbf{x} \in (\mathcal{S} \times \mathcal{A} \times \mathbb{R}_0)^{nM} \mid \Phi^\pi = \Gamma^{\pi^\dagger}\}$ に対して，

$$\Pr(\Phi^\pi = \Gamma^{\pi^\dagger}) \approx \frac{|\mathcal{X}_n^\dagger|}{\exp\{n(H(\Gamma^\pi | \mathbf{V}) + H(\Gamma^\pi | \mathbf{W}) + \epsilon_n)\}}, \quad (26)$$

となる．ただし， ϵ_n は $\epsilon_n \rightarrow 0$ ($n \rightarrow \infty$) となる数である．MAS を MDP で定常エルゴードと見なせるほどゆっくりと学習させた場合，例えば最適政策行列 Γ^{π^\dagger} の各要素が決定論的 (0 もしくは 1) ならば，明らかに $H(\Gamma^\pi | \mathbf{V}) \rightarrow 0$ となる．その結果，(26) において，収益最大化の確率の指数的増加の過程が確認できる．

4 まとめ

MAS が定常エルゴード MDP に従う場合に成立する漸近的性質 (22)–(24) を示し，最適政策の学習により MAS の収益が最大化される過程を (26) により明らかにした．

参考文献

- [1] G. Chen *et al.*, Coordinating Multiple Agents via Reinforcement Learning, *Autonomous Agents and Multi-Agent Systems*, vol. 10, no. 3, pp. 273–328, 2005.
- [2] H. J. Kushner and G. G. Yin, *Stochastic Approximation Algorithms and Applications*, Springer-Verlag, 1997.
- [3] I. Csiszár, The Method of Types, *IEEE Trans. Information Theory*, vol. 44, no. 6, pp. 2505–2523, 1998.
- [4] K. Iwata *et al.*, A Statistical Property of Multiagent Learning Based on Markov Decision Process, *IEEE Trans. Neural Networks*, vol. 17, no. 4, pp. 829–842, 2006.