

マルコフ型決定過程(I)

日大 生産工 坂本 武司

1. 平均最適政策

有限個の状態、決定からなる離散型マルコフ型決定過程に
ついて述べる。

S 位の状態 $1, 2, \dots, S$ からなる系 (System) を考之。離
散時刻 ($n = 1, 2, \dots$) ごとに状態を観測し、可能な決定の
集合 $K = \{k\} = \{1, 2, \dots, K\}$ から決定を選択。もし、 n 時
刻の状態 i を観測し、決定 k を行えば、その結果

- (1) 系は推移確率 q_{ij}^k によって新しい状態 j に移る。
- (2) 状態 i で決定 k を行えば、利得 r_i^k を得る。

記号.

状態空間: $S = \{1, 2, \dots, S\}$

決定空間: $K = \{1, 2, \dots, K\}$

推移確率: q_{ij}^k , $\sum_{j=1}^S q_{ij}^k = 1$

利得: r_i^k

f を S から K への関数とし、 f の全体を F と表す。下の要素の系列 $\pi = (f_1, f_2, \dots, f_n, \dots)$ を政策と呼ぶ。政策 π を用いるとは初期状態が i_1 あれば、第1期で決定 $f_1(i_1)$ 、第2期の状態が i_2 あれば $f_2(i_2)$ 以下同様に続けて、第n期で $f_n(i_n)$ を行う。 $\pi = (f, f, \dots, f, \dots) \equiv f^\infty$ を定常政策と呼ぶ。

$f \in F$ に対する利得及び確率ベクトルを次のように表す。

$$r(f) = (r_1^{f(1)}, r_2^{f(2)}, \dots, r_s^{f(s)})' \quad (S \times 1 \text{ ベクトル})$$

$$Q(f) = (q_{ij}^{f(i)}) \quad (S \times S \text{ マトリクス})$$

$V_n(i, \pi)$ = 状態 i から出発して、政策 π を用いると次の n 期の総期待利得とす。このベクトルを

$$V_n(\pi) = (V_n(1, \pi), \dots, V_n(S, \pi))' \quad (S \times 1 \text{ ベクトル})$$

とする。そのとき

$$(1.1) \quad V_n(\pi) = \sum_{k=0}^{n-1} Q_k(\pi) r(f_{n+k})$$

但し、 $Q_k(\pi) = Q(f_1) \cdots Q(f_k)$, $Q_0(\pi) = I$ (単位行列)

注意。 $S \times 1$ ベクトル w に対して作用素 $L(f)$ を

$$(1.2) \quad L(f)w = r(f) + Q(f)w$$

と定義する。 $L(f)$ は monotone である。即ち、注意。ベクトル u, w に対して、すべての i で $u_i \geq w_i$ ならば $u \geq w$ 且 $u \neq w$ ならば $u > w$ と定義すると、 $u \geq w$ 且

らば、 $L(f)U \geq L(f)W$ である。今後上記のよろべストルの順序を定めておく。 $L(f)$ を用いると。

$$V_n(\pi) = L(f_1)L(f_2)\cdots L(f_{n-1})r(f_n)$$

$V_n(\pi)$ は $n \rightarrow \infty$ のとき発散し、1期当たりの平均利得 $\frac{1}{n} V_n(\pi)$ は一般に収束して $\underline{U}(f)$ である。 $\underline{U}(f) = \liminf_{n \rightarrow \infty} \frac{1}{n} V_n(\pi)$ を考えると $U(f) \geq \underline{U}(f)$ for all π を満たす定常政策 f^* が存在する。但し $U(f) = \underline{U}(f) = Q^*(f)r(f)$ 。また割引率 β を入れると、 $\beta \rightarrow 1$ のとき、無限過程の終期待利得を最大にする定常政策が存在する。更に、各 $n \geq n_2$ 、 $U(n) = \max_{\pi} V_n(\pi)$ とするとき $\lim_{n \rightarrow \infty} \frac{1}{n} U(n) = U(f)$ となる定常政策 f^* が存在する。そして、この3つの定常政策の平均利得は一致する。そこで、定常政策の中で、平均利得を最大にする政策を決定する algorithm を逐一換及び割引率 β を用いて求めよう。

Theorem 1. (1), (2). 任意の $f \in \mathbb{F}$ は \mathbb{Z} 、 $Q^n(f) \rightarrow Q^*(f)$ ($n \rightarrow \infty$) とする。

$$(1.1) \quad V_n(f^*) = nU(f) + v(f) + \varepsilon(n, f)$$

(b). $E(f)$ のエルゴディック集合の周期の最小公倍数を N とし、

$$Q_0 = Q^N$$

$$(1.2) \quad V_{nN+m}(f^*) = (nN+m)U(f) + v(f) + w(m, f) + \varepsilon(n, f)$$

従し、 $U(f)$ は

$$(1.3) \quad (I - Q(f))u = 0 \quad Q^*(f)u = Q^*(f)r(f)$$

の一意の解であり、 $v(f)$ は

$$(1.5) \quad (I - Q(f)) v = r(f) - u(f), \quad Q^*(f)v = 0$$

の一意の解である。また、 $w(nN+m, f) = w(m, f) \quad (n=1, 2, \dots)$
 $\varepsilon(n, f) \rightarrow 0 \quad (n \rightarrow \infty)$

$$(2) \quad G_1(f) \equiv \{ g \mid g \in F, \quad Q(g)u(f) > u(f) \}$$

$$\begin{aligned} G_2(f) \equiv \{ g \mid g \in F, \quad Q(g)u(f) &= u(f), \quad r(g) + Q(g)v(f) \\ &> u(f) + v(f) \} \end{aligned}$$

$$G(f) \equiv G_1(f) \cup G_2(f) \text{ とす。}$$

$g \in G(f)$ ならば、 $u(g) \geq u(f)$

$$(3) \quad G(f) = \emptyset \text{ ならば、任意の } g \in F \Rightarrow n \text{ で } u(f) \geq u(g)$$

$$\begin{aligned} (\text{証明}) \quad (1) \quad (2). \quad V_m(f^\infty) &= \sum_{k=0}^{n-1} Q^k(f) r(f) \\ &= (m-1) Q^*(f) r(f) + \sum_{k=0}^{m-1} (Q(f) - Q^*(f))^k r(f) \end{aligned}$$

$n \rightarrow \infty$ のとき、 $Q^n(f) \rightarrow Q^*(f)$ となる。 n 分+余分でないとき。

$$\begin{aligned} V_m(f^\infty) &= n Q^*(f) r(f) + H(f) r(f) + \varepsilon(n, f) \\ &= n u(f) + v(f) + \varepsilon(n, f) \end{aligned}$$

但し、 $v(f) \equiv H(f) r(f)$, $H(f) \equiv (I - Q(f) + Q^*(f))^{-1} - Q^*(f)$

$$\varepsilon(n, f) \rightarrow 0 \quad (n \rightarrow \infty)$$

$$(b) \quad Q^*(f) = \frac{1}{N} Q_0^*(f) \sum_0^{N-1} Q^i(f) \text{ より}$$

$$\begin{aligned} V_{nN+m}(f^\infty) &= (m-1) N Q^*(f) r(f) + [I - (Q_0(f) - Q_0^*(f))]^{-1} \\ &\quad \sum_0^{N-1} Q^i(f) r(f) + Q_0^*(f) \sum_0^{m-1} Q^i(f) r(f) \end{aligned}$$

よって、2点より上の結果を得る。

(2). 周期的では場合 (2) の 2 示す。仮定より n を十分大きくて

$$\begin{aligned} V_{nN+m+1}(g, f^\infty) &= Q(g) W(m, f) + Q(g) \varepsilon(nN+m, f) \\ &> V_{nN+m+1}(f^\infty) = W(m+1, f) + \varepsilon(nN+m+1, f) \end{aligned}$$

$$\text{即ち}, L(g) V_{nN+m}(f^\infty) > V_{nN+m+1}(f^\infty) + Q(g) W(m, f) \\ - \varepsilon(nN+m+1, f) + Q(g) \varepsilon(nN+m, f)$$

一般に、往々 g の各整数 M について。

$$\begin{aligned} L^M(g) V_{nN+m}(f^\infty) &> V_{nN+m+M}(f^\infty) + Q^M(g) W(m, f) \\ &- W(m+M, f) - \varepsilon(nN+m+M, f) + Q^M(g) \varepsilon(nN+m, f) \end{aligned}$$

実際、 $M = 1$ のとき明らかに $L(g) V_{nN+m}(f^\infty) > V_{nN+m+1}(f^\infty) + Q(g) W(m, f) - Q(g) W(m+1, f) - Q(g) \varepsilon(nN+m+1, f) + Q(g) \varepsilon(nN+m, f)$

$$\begin{aligned} L^{M+1}(g) V_{nN+m}(f^\infty) &\geq L(g) V_{nN+m+M}(f^\infty) + Q^{M+1}(g) W(m, f) \\ &- Q(g) W(m+M, f) - Q(g) \varepsilon(nN+m+M, f) \\ &+ Q^{M+1}(g) \varepsilon(nN+m, f) \end{aligned}$$

仮定を用いて。

$$\begin{aligned} L^{M+1}(g) V_{nN+m}(f^\infty) &> V_{nN+m+1}(f^\infty) + Q^{M+1}(g) W(m, f) \\ &- W(m+M+1, f) - \varepsilon(nN+m+M+1, f) + Q^{M+1}(g) \varepsilon(nN+m, f) \end{aligned}$$

$$\min_{\substack{1 \leq i \leq s \\ 1 \leq m \leq N}} w_i(m, f) = c_1, \quad \max_{\substack{1 \leq i \leq s \\ 1 \leq m \leq N}} w_i(m, f) = c_2, \quad \delta = (1, b, \dots, 1)$$

とおこう。

$$\begin{aligned} L^M(g) V_{nN+m}(f^\infty) &> V_{nN+m+M}(f^\infty) + (c_1 - c_2) \delta \\ &- \varepsilon(nN+m+M, f) + Q^M(g) \varepsilon(nN+m, f) \end{aligned}$$

$M \rightarrow \infty$ とすると、 $\varepsilon(nN+m+M, f) \rightarrow 0$ また $\varepsilon(nN+m, f)$

Q^M は有界だから、上式を M を割り $M \rightarrow \infty$ とすると、 $U(g) \geq U(f)$ を得る。周期的でないときも同様である。

3.2. 強化最適政策の計算法

割引率 β ($0 \leq \beta < 1$) を考えたときの政策 π に対する終期待利得を $V_\beta(\pi)$ と表すと、

$$(2.1) \quad V_\beta(\pi) = \sum_{n=0}^{\infty} \beta^n Q_n(\pi) r(f_{n+1}) \quad \text{但し } Q_0(\pi) = I$$

もし、 $U(\beta) \equiv V_\beta(\pi^*) \geq V_\beta(\pi)$ for all π が成立すれば、 π^* を β -最適政策と呼ぶ。

また、適当な β_0 が存在して ($0 \leq \beta_0 < 1$)、すべての $\beta_0 \leq \beta < 1$ に対して $V_\beta(\pi^*) \geq V_\beta(\pi)$ for all π が成立すれば、 π^* を最適政策と呼ぶ。 β -最適又は最適な定常政策が存在するとは先づ知らねえ。

さて、 $U(\beta) - V_\beta(\pi) \rightarrow 0$ ($\beta \uparrow 1$) のときは π を強化最適 (nearly optimal) 又は 1-最適 (1-optimal) と呼ぶ。これは、強化最適政策を決定する Veinott's algorithm 12) を述べる。

まず、準備として D. Blackwell の結果を述べる。

Theorem 1 と同じ手順により

Lemma 1. 任意の $f \in F$ は \exists

$$(2.2) \quad V_\beta(f^\infty) = \frac{u(f)}{1-\beta} + v(f) + e(\beta, f), \quad 0 \leq \beta < 1$$

$$\varepsilon(\beta, f) \rightarrow 0 \quad (\beta \rightarrow 1)$$

次の集合を定義する。

$$F' = \{f \mid f \in F, u(f) \geq u(g) \text{ all } g \in F\}$$

$$F'' = \{f \mid f \in F', v(f) \geq v(g) \text{ all } g \in F'\}$$

(2.2) 算法より

Lemma 2. F'' は殆ど最適なすべての $f \in F$ の集合である。

1. たとえば 2. F'' の要素を決定する algorithm を作ればいい。

。定理 1 より $G(f) = \emptyset$ なら $f \in F'$ であるから、 F' の中で v を最大にする計算法を求めるばよ。

Lemma 3. $f \in F, g \in G(f)$ ならば: $u(g) > u(f)$ 又は:

$$u(g) = u(f), v(g) > v(f)$$

(証明). 定理 1 より、 $g \in G(f)$ ならば $u(g) \geq u(f)$. 又: $u(g) = \lim_{n \rightarrow \infty} \frac{1}{n+1} \sum_{k=0}^n Q^k(f) r(f) = \lim_{\beta \rightarrow 1} (1-\beta) \sum_{k=0}^{\infty} \beta^k Q^k(f) r(f) = \lim_{\beta \rightarrow 1} (1-\beta) V_\beta(f^\infty)$ す。 β が 1 に近づくとき $V_\beta(f^\infty) > V_\beta(f)$

$$V_\beta(g^\infty) = \frac{u(g)}{1-\beta} + v(g) + \varepsilon(\beta, g)$$

$$V_\beta(f^\infty) = \frac{u(f)}{1-\beta} + v(f) + \varepsilon(\beta, f)$$

である。 $u(g) \geq u(f)$ 且つ $u(g) = u(f)$ ならば $v(g) \geq v(f)$

もし $(u(g), v(g)) = (u(f), v(f))$ ならば $g \in G(f)$ は矛盾。

次の集合を定義する。

$$E(f) = \{g \mid g \in F, Q(g)u(g) = u(f), n(g)Q(g)v(f) = v(f)\}$$

$E(f)$ の定義と (1.3), (1.4) より $g \in E(f)$ ならば $u(g) = u(f) < v(f) < v(g)$

$\Gamma \vdash f \rightarrow z$ は $f \in \Gamma' \vdash z$ は $E(f) \subset \Gamma'$ である。集合 $E(f)$ は計算上
丰子から $E(f)$ の中で V を最大にするものを考へる。

Lemma 4. (Veinott) $f \in F$, $g \in E(f)$ とする。 $w(g)$ は
(2.3) $[I - Q(g)] w = 0$, $Q^*(g)w = Q^*(g)(-v(f))$
の一意の解とする。証明を

$$(2.4) \quad v(g) = v(f) + w(g)$$

(証明) 一意性は定理 2 より出る。 $g \in E(f)$ より

$$[I - Q(g)] v(f) = r(g) - u(f)$$

この式と (2.3) を組み合わせて、

$$[I - Q(g)][v(f) + w(g)] = r(g) - u(f), \quad Q^*(g)[v(f) + w(g)] = 0$$

よって定理 2 より、(2.4) が成立する。

この Lemma と $E(f)$ の中で $w(g)$ を最大にすることは、 F を $E(f)$ で $r(f)$ を $-v(f)$ で置きかえれば、 F の中で $u(f)$ を最大にするのと同じ形をしていふことわかる。従って $E(f)$ の中で $v(g)$ を最大にするためには、政策反復法を再び使之はよい。
ことを Γ とする。

任意の $f \in F$ について $\Xi(f)$ を (2.5) の一意の解とする。

$$(2.5) \quad (I - Q(f)) \Xi(f) = -v(f), \quad Q^*(f) \Xi(f) = 0$$

$$(2.6) \quad H(f) \equiv \{ g \mid g \in E(f), -v(f) + Q(g) \Xi(f) > \Xi(f) \}.$$

とする。証明を

Theorem 2. (Veinott)

(1) $f \in F$, $G(f) = \phi$ とする. すなはち $V(f) \geq V(g)$ for all

$g \in E(f)$ ならば $f \in F'$

(2) $G(f) \cup H(f) = \phi$ ならば $f \in F'$

(3) $g \in H(f)$ ならば $U(g) = U(f)$ であり、更に $V(g) > V(f)$

又は $V(g) = V(f)$, $Z(g) > Z(f)$ のいずれかが成立する。

(証明) (1) $g \in F' - E(f)$ とする. $G(f) = \phi$, $g \in F'$ より $U(g)$

$= U(f)$. 更に $g \notin E(f)$ より $r(g) + Q(g)V(f) < U(f) + V(f)$

$$V_\beta(g, f^\infty) = \frac{Q(g)U(f)}{1-\beta} + r(g) - Q(g)U(f) + Q(g)V(f) + \varepsilon(\beta, f, g)$$

$$V_\beta(f^\infty) = \frac{U(f)}{1-\beta} + V(f) + \varepsilon(\beta, f)$$

であるから、十分 $1 < \beta < 1 - \varepsilon$. $V_\beta(g, f^\infty) < V_\beta(f^\infty)$

よって $V_\beta(g^\infty) < V_\beta(f^\infty)$ (1) 参照)、即ち $V(f) > V(g)$.

(2) $H(f) = \phi$ であるから、 $W(f) \geq W(g)$ for all $g \in E(f)$. また

又は $W(f) = 0$ であるから、 $V(g) = V(f) + W(g) \leq V(f) + W(f) = V(f)$

for all $g \in E(f)$.

(3) $g \in H(f) \subset E(f)$ より $U(g) = U(f)$. 更に Lemma 3 より

$W(g) > W(f) = 0$. 又は $W(g) = W(f) = 0$, $Z(g) > Z(f)$.

この定理より、 $f_1 \in F$ とし、 f_2, f_3, \dots を $f_{i+1} \in G(f_i) \cup H(f_i)$ であるように選ぶと、 $\{U_i(f), V_i(f), Z_i(f)\}$ (は辞書的順序で増加するから、同じ f_i が 2 度現われることはない)。F は有限集合だから、ある i で $G(f_i) \cup H(f_i) = \phi$ となり、 f_i が最適政策を得られる。

例. $S = \{1, 2\}$, $K_1 = \{1, 2\}$, $K_2 = \{1\}$ (K_i は状態 i の可能行決定の集合)

$$F = \{f, g\} \text{ 且し } f(i) = 1 \quad (i=1, 2)$$

$$f(1) = 2, \quad g(2) = 1$$

$$r(f) = \begin{pmatrix} 3 \\ -3 \end{pmatrix} \quad Q(f) = \begin{pmatrix} \frac{1}{2} & \frac{1}{2} \\ \frac{1}{2} & \frac{1}{2} \end{pmatrix} \quad Q^*(f) = \begin{pmatrix} \frac{1}{2} & \frac{1}{2} \\ \frac{1}{2} & \frac{1}{2} \end{pmatrix}$$

$$r(g) = \begin{pmatrix} 6 \\ -3 \end{pmatrix} \quad Q(g) = \begin{pmatrix} 0 & 1 \\ \frac{1}{2} & \frac{1}{2} \end{pmatrix} \quad Q^*(g) = \begin{pmatrix} \frac{1}{3} & \frac{2}{3} \\ \frac{1}{3} & \frac{2}{3} \end{pmatrix}$$

$$u(f) = \begin{pmatrix} 0 \\ 0 \end{pmatrix} \quad v(f) = \begin{pmatrix} 3 \\ -3 \end{pmatrix} \quad z(f) = \begin{pmatrix} -3 \\ 3 \end{pmatrix}$$

$$u(g) = \begin{pmatrix} 0 \\ 0 \end{pmatrix} \quad v(g) = \begin{pmatrix} 4 \\ -2 \end{pmatrix} \quad z(g) = \begin{pmatrix} -\frac{8}{3} \\ \frac{4}{3} \end{pmatrix}$$

$$G(f) = \emptyset, \quad G(g) = \emptyset, \quad H(f) = g, \quad H(g) = \emptyset$$

g が3台比最適政策

参考文献

- (1) D. Blackwell, Discrete Dynamic Programming, Ann. Math. Statist. 33 (1962) 719-726
- (2) B. W. Brown, On the Iterative Method of Dynamic Programming on a Finite Space Discrete Time Markov Processes, Ann. Math. Statist. 36 (1965) 1279-1285
- (3) C. Deiman, On Sequential Decisions and Markov

chains, Management Science 9 (1962) 16-24

(4) R.A. Howard, Dynamic Programming and Markov Processes, 1960, Technology Press and Wiley

(5) M. Ogawara, A Note on Discrete Markovian Decision Processes, Bull. Math. Statist 11 (1963)

35-42

(6) A.F. Veinott, Jr., On Finding Optimal Policies in Discrete Dynamic Programming, Ann. Math. Statist 37 (1966) 1284-1294