

## 多重連鎖マルコフ決定過程の逐次近似法について

名古屋工業大学 大野 勝久 (Katsuaki Ohno)

### 1. 序論

有限状態、有限決定をもつマルコフ決定過程について、時間平均利得を最大にする最適定常政策を決定するアルゴリズムとしては、政策反復法、逐次近似法、線形計画法、修正政策反復法が知られている。特に逐次近似法、修正政策反復法は多状態問題に対する有力な手法として多くの研究が行われておらず、多重連鎖問題については余り論じられていない。

Bather (1973) は多重連鎖問題にたいする communicating sets への分解と政策反復法に関する逐次近似法を提案し、その収束を示している。また  $\pi$ -最適政策の構成をも論じているが、反復回数が十分大きければ  $\pi$ -最適政策がえられると言べるにとどまっている。Schweitzer (1984) もまた communicating sets への分解 (unique chain decomposition と名づけている) にもと

づく逐次近似アルゴリズムを提案し、その収束を示していふが、 $\varepsilon$ -最適政策を保証する停止基準は与えられていない。本論文では  $\varepsilon$ -最適政策を求める逐次近似法を Bather (1973) を参考に論ずる。

## 2. 準備

以下の記号を使用する。

$I = \{1, 2, \dots, M\}$  : 状態空間

$K_i (i \in I)$  : 状態  $i$  でヒリうる決定の有限集合

$r_i(k) (i \in I, k \in K_i)$  :  $i$  において  $k$  をとったときの平均利得

$p_{ij}(k) (i, j \in I, k \in K_i)$  :  $i$  において  $k$  をとった時に遷移する確率

$F$  : 定常政策  $f = (f_1, \dots, f_M)$  の集合 ( $f_i \in K_i, i \in I$ )

$r(f) = (r_1(f_1), \dots, r_M(f_M))^T$  (T は転置を表す)

$P(f) = (p_{ij}(f_i))$  :  $f$  のもとでの遷移確率行列

$P^*(f)$  :  $P(f)$  のチエザロ和 ( $\lim_{N \rightarrow \infty} \frac{1}{N+1} \sum_{n=0}^N P(f)^n$ )

$g(f) = P^*(f)r(f)$  :  $f$  の利得 (gain)

$v(f)$  :  $f$  の相対値

$g(f)$  は

$$(1) \quad g = P(f)g$$

$$(2) \quad g + v = r(f) + P(f)v$$

の一意解であり、 $v(f)$  は  $P(f)$  の各エルゴード連鎖に属する

この状態  $i$  で  $v(f) = 0$  とおけば一意に決定される。時間平均マルコフ決定過程は 4 つ組  $(I, r, P, F)$  で与えられ、利得を最大にする政策  $f^* \in F$  を決定する問題である。

最大利得  $g^* = g(f^*) = \max_{f \in F} P^*(f) r(f)$  および相対値  $v^* = v(f^*)$  のみたす最適方程式は

$$(3) \quad g_i^* = \max_{k \in K_i} \left\{ \sum_{j \in I} p_{ij}(k) g_j^* \right\} \quad (i \in I)$$

$$(4) \quad g_i^* + v_i^* = \max_{k \in L_i} \left\{ r_i(k) + \sum_{j \in I} p_{ij}(k) v_j^* \right\} \quad (i \in I)$$

である。ここで

$$L_i = \{ k \in K_i ; (3) 式右辺を最大にする k \}$$

である。定常政策子は、 $e = (1, \dots, 1)^T$  に対して

$$(5) \quad g^* - g(f) \leq \varepsilon e$$

をみたすとき、 $\varepsilon$ -最適政策と呼ばれる。

### [定理 1]

$f \in F$  および  $M$  次元ベクトル  $g, v$  にたいして

$$(6) \quad \psi = g - P(f)g, \quad \gamma = g + v - r(f) - P(f)v$$

とおく。  $P(f)$  が

$$(7) \quad P(f) = \begin{pmatrix} Q(f) & 0 \\ R(f) & S(f) \end{pmatrix} \quad (S(f) \text{ は過渡行列})$$

とすると、 $Q(f), S(f)$  に対応する部分ベクトルを添字  $c, t$  で表わしたとき

$$(8) \quad \psi_c \leq (\geq) 0, \quad \psi_t \leq (\geq) a_t, \quad \gamma_c \leq (\geq) b_c$$

が" たとえ T = T<sub>c</sub> ば。

$$(9) \quad g_c - g_c(f) \leq (\geq) Q^*(f) b_c$$

$$(10) \quad g_t - g_t(f) \leq (\geq) (I - S(f))^{-1} (a_t + R(f) Q^*(f) b_c)$$

である。すなはち I は単位行列を表わす。

(証明)  $\Delta g = g - g(f)$ ,  $\Delta v = v - v(f)$  とおけば、(1), (2) 式より

$$(11) \quad \Delta g = \psi + P(f) \Delta g$$

$$(12) \quad \Delta g + \Delta v = \gamma + P(f) \Delta v$$

である。また、(8), (11) 式より  $\Delta g_c \leq (\geq) Q^*(f) \Delta g_c$  である。すなはち、(8), (12) 式より

$$\Delta g_c + \Delta v_c \leq (\geq) b_c + Q(f) \Delta v_c$$

である、 $Q^*(f)$  を左からかけねば  $Q^*(f) Q(f) = Q^*(f)$  であるから。

(9) 式よりえられる。また (8), (11) 式より

$$\Delta g_t \leq (\geq) a_t + R(f) \Delta g_c + S(f) \Delta g_t$$

である、 $S(f)$  は過渡行列であるから (9) 式より (10) 式がえられる。

定理 1 において、 $g = g^*$ ,  $v = v^*$ ,  $a_t = \varepsilon_1 e_t$ ,  $b_c = \varepsilon_2 e_c$  とおけば、  
 $(I - S(f))^{-1} R(f) e_c = e_t$  すなはち次の系がえられる。

[系 1]

(7) 式でえられる  $P(f)$  について

$$(13) \quad g_c^* \leq Q(f) g_c^*, \quad g_t^* - (R(f) g_c^* + S(f) g_t^*) \leq \varepsilon_1 e_t$$

$$g_c^* + v_c^* - R(f) - Q(f) v_c^* \leq \varepsilon_2 e_c$$

が" たてば"

$$(14) \quad g_c^* - g_c(f) \leq \varepsilon_2 e_c, \quad g_t^* - g_t(f) \leq \varepsilon_2 e_t + \varepsilon_1 (I - S(f))^{-1} e_t$$

である。

Whittle (1983, 定理 4.1(b), p.122) は最小化問題について、系 1 と同様  $f^*$  が  $\varepsilon$ -最適となる条件を示しているが、証明中にミスゴ"あり、示された条件は正しくない。

定理 1 の ( ) に対応する式において  $f = f^*$  とおけば、系 1 と同様にて次の系せえられる。

[系 2]

$P(f^*)$  が (7) 式の形で与えられるとき、

$$(15) \quad g_c \geq Q(f^*) g_c, \quad g_t - (R(f^*) g_c + S(f^*) g_t) \geq -\varepsilon_1 e_t \\ g_c + v_c - r_c(f^*) - Q(f^*) v_c \geq -\varepsilon_2 e_c$$

が" たてば"

$$(16) \quad g_c^* - g_c \leq \varepsilon_2 e_c, \quad g_t^* - g_t \leq \varepsilon_2 e_t + \varepsilon_1 (I - S(f^*))^{-1} e_t$$

である。

### 3. 最適停止問題

Bather (1973), Federgrun and Schweitzer (1984), Schweitzer (1984) に従ふ。まず状態空間  $I$  を communicating sets へ分解する。

(17) Unique Chain Decomposition :

1.  $I_1 = I$ ,  $D_i = K_i$  ( $i \in I_1$ ),  $\lambda = 1$  とおく。

2. 各  $i \in I_\ell$  において  $D_i$  の全ての決定を正の確率でとる政策

にたいする  $I_l$  の全ての部分連鎖  $C(r; l)$  ( $r = 1, \dots, R_l$ ) を定めよ。

3.  $T_l = \phi$ ,  $J = I_l - \sum_r C(r; l)$  とおく。 $i \in J$  にたいして  $\sum_{j \in J} p_{ij}(k) < 1$  となる  $k$  を  $D_i$  からとり除く,  $D_i = \phi$  となれば  $i$  を  $J$  から  $T_l$  に移す。この操作を  $D_i$  ( $i \in J$ ) が変化しなくなるまで反復す

る。

4.  $J = \phi$  となれば  $L = l$  とおいて 5. へ。ともなければ、

$I_{l+1} = J$ ,  $l = l+1$  とおいて 2. へ。

5.  $R = \sum_{l=1}^L R_l$ ,  $T = \sum_{l=1}^L T_l$  とおく,  $l = 1, \dots, L$ ,  $r = 1, \dots, R_l$  にたいして  $C\left(\sum_{m=1}^{l-1} R_m + r\right) = C(r; l)$  とおく。

上記アルゴリズムからえられる各  $C(r)$  ( $r \in R = \{1, \dots, R\}$ ) は、  
 $D_i$  ( $i \in C(r)$ ) からの政策にたいして communicating set をなす。定  
常政策  $f(r) = (f_i)$  ( $f_i \in D_i$ ,  $i \in C(r)$ ) の集合を  $F(r)$  とおけば、部分マルコフ決定過程  $(C(r), r, P, F(r))$  ( $r \in R$ ) えられる。  
 $(C(r), r, P, F(r))$  は状態に依存しない最大利得  $\sigma^*(r)$  をもち、  
相対値  $u_i^*(r)$  ( $i \in C(r)$ ) は最適方程式

$$(18) \quad \sigma^*(r) + u_i^*(r) = \max_{k \in D_i} \left\{ r_i(k) + \sum_{j \in C(r)} p_{ij}(k) u_j^*(r) \right\} \quad (i \in C(r))$$

の解である。

(18) 式を満たす  $\sigma^*(r)$ ,  $u_i^*(r)$  および最適政策  $f^*(r)$  ( $r \in R$ ) は  
いてえられるものとする。このとき  $i \in C(r)$  にたいして  $\tilde{K}_i =$   
 $K_i - D_i$  となれば (3) または

$$(19) \quad g_i^* = \max \left\{ \sigma^*(r), \max \left\{ \sum_{j \in I} p_{ij}(k) g_j^* \mid k \in \tilde{K}_i \right\} \right\}$$

とある。アルゴリズム (17) によると  $r \leq R_1$  については  $\tilde{K}_i = \emptyset$   
( $i \in C(r)$ ) であるから

$$(20) \quad g_i^* = \sigma^*(r) \quad (i \in C(r), r \leq R_1)$$

である。 $r > R_1$  については  $C(r)$  の下にとも 1 つの  $i$  で  $\tilde{K}_i \neq \emptyset$  である, すな

$$(21) \quad \sigma^*(r) > \max_{i \in C(r)} \max_{k \in \tilde{K}_i} \left\{ \sum_{j \in I} P_{ij}(k) g_j^* \right\} = \sum_{j \in I} P_{i^* j}(k^*) g_j^*$$

となるれば  $i \in C(r)$  で (20) 式が成立する。一方,

$$(22) \quad \sum_{j \in I} P_{i^* j}(k^*) g_j^* \geq \sigma^*(r)$$

となるば、状態  $i^*$  で決定  $k^* \in \tilde{K}_{i^*}$  をとり、 $i^*$  以外の  $i \in C(r)$  では  $i^*$  を吸収状態 とするような政策  $f(r) \in F(r)$  をとれば全ての  $i \in C(r)$  で利得  $\sum_{j \in I} P_{ij}(k^*) g_j^*$  がえられる。ゆえに  $i \in C(r)$  ( $r \in R$ ) については

$$(23) \quad g^*(r) = g_{i^*}^*, \quad \tilde{K}(r) = \{(i, k) ; i \in C(r), k \in \tilde{K}_i\}$$

である。 $\hat{k} = (i, k) \in \tilde{K}(r)$  については

$$(24) \quad \hat{P}_{rs}(\hat{k}) = \sum_{j \in C(s)} P_{ij}(k) \quad (s \in R), \quad \hat{P}_{rj}(\hat{k}) = P_{ij}(k) \quad (j \in T)$$

である。他方  $i \in T$  については  $k \in K_i$  については

$$(25) \quad \hat{P}_{is}(k) = \sum_{j \in C(s)} P_{ij}(k) \quad (s \in R), \quad \hat{P}_{ij}(k) = P_{ij}(k) \quad (j \in T)$$

であることに注意。こゝと同様の最適停止問題が導かれる。

$$(26) \quad \hat{g}_i^*(r) = \max \left\{ \sigma^*(r), \max \left\{ \sum_{s=1}^R \hat{P}_{rs}(\hat{k}) \hat{g}_s^*(s) + \sum_{j \in T} \hat{P}_{rj}(\hat{k}) \hat{g}_j^* \mid \hat{k} \in \tilde{K}(r) \right\} \right\} \quad (r \in R)$$

$$\hat{g}_i^* = \max \left\{ \sum_{s=1}^R \hat{P}_{is}(k) \hat{g}_s^*(s) + \sum_{j \in T} \hat{P}_{ij}(k) \hat{g}_j^* \mid k \in K_i \right\} \quad (i \in T)$$

∴  $\hat{K}(r) = \emptyset$  ( $r \leq R_1$ ) である。各  $r \in R$  において決定 "停止"

(以後 "0" で表わす) を選べば、利得  $\sigma^*(r)$  をえて過程は停止する。問題 (26) にいたる正常政策  $\tilde{f}$  は、決定  $\tilde{f}(r) \in \tilde{K}(r) \cup \{0\}$  ( $r \in R$ ),  $\tilde{f}_i \in K_i$  ( $i \in T$ ) の組  $(\tilde{f}(r), \tilde{f}_i)$  で与えられる。 $\tilde{f}$  に  $T = \{\cdot\}$  とすると、 $R_c(\tilde{f}) = R - R_s(\tilde{f})$  とかく。全ての正常政策  $\tilde{f}$  の集合を  $\tilde{F}$ 、問題 (26) の最適政策を  $\tilde{f}^*$  で表わす。

$R + |T|$  次ベクトル  $\tilde{g} \in r \in R, \tilde{g} \in \tilde{K}(r), i \in T, k \in K_i$  にて次の写像  $T, U$  を

$$(27) \quad T(r; \tilde{g}) \tilde{g} = \sum_{s \in R} \tilde{P}_{rs}(\tilde{g}) \tilde{g}(s) + \sum_{j \in T} \tilde{P}_{rj}(\tilde{g}) \tilde{g}_j$$

$$T(i; k) \tilde{g} = \sum_{s \in R} \tilde{P}_{is}(k) \tilde{g}(s) + \sum_{j \in T} \tilde{P}_{ij}(k) \tilde{g}_j$$

$$(28) \quad U(r) \tilde{g} = \max \{ T(r; \tilde{g}) \tilde{g} \mid \tilde{g} \in \tilde{K}(r) \}$$

$$U(i) \tilde{g} = \max \{ T(i; k) \tilde{g} \mid k \in K_i \}$$

で定義する。政策  $\tilde{f}$  へもとての利得を  $\tilde{g}(\tilde{f}) = (\tilde{g}(r; \tilde{f}), \tilde{g}(i; \tilde{f}))$  で表わせば、 $\tilde{g}(\tilde{f})$  は

$$(29) \quad \tilde{g}(r; \tilde{f}) = \sigma^*(r) \quad (r \in R_s(\tilde{f}))$$

$$(30) \quad \tilde{g}(r; \tilde{f}) = T(r; \tilde{f}(r)) \tilde{g}(\tilde{f}) \quad (r \in R_c(\tilde{f}))$$

$$\tilde{g}(i; \tilde{f}) = T(i; \tilde{f}_i) \tilde{g}(\tilde{f}) \quad (i \in T)$$

を満たしていき。 $r \in R_s(\tilde{f})$  にいたりベクトル  $\tilde{g}_s(\tilde{f}) = (\tilde{g}(r; \tilde{f}),$   
 $\sigma^* = (\sigma^*(r))$  を導入し、 $r \in R_c(\tilde{f})$  および  $i \in T$  にいたりベクトル  
 $\tilde{g}_{ct}(\tilde{f}) = (\tilde{g}(r; \tilde{f}), \tilde{g}(i; \tilde{f}))^T$ ,  $T_{ct}(\tilde{f}) \tilde{g} = (T(r; \tilde{f}(r)) \tilde{g}, T(i; \tilde{f}_i) \tilde{g})^T$  を導入すれば、(29), (30) 式は各々

$$(29)' \quad \tilde{g}_s(\tilde{f}) = \sigma_s^*$$

$$(30)' \quad \tilde{g}_{ct}(\tilde{f}) = T_{ct}(\tilde{f}) \tilde{g}(\tilde{f})$$

と書き直さぬ。

### [補題 1]

任意の  $\tilde{f} \in \tilde{F}$  にてして,  $Rc(\tilde{f})UT$  から  $Rc(f)UT$  への遷移行列

$$(31) \quad \tilde{\delta}(\tilde{f}) = \begin{pmatrix} (\tilde{P}_{rs}(\tilde{f}(r))) & (\tilde{P}_{rj}(\tilde{f}(r))) \\ (\tilde{P}_{is}(\tilde{f}_i)) & (\tilde{P}_{ij}(\tilde{f}_i)) \end{pmatrix}$$

は過渡行列である。 $\tilde{g}(\tilde{f})$  は  $r \in R(\tilde{f})$  にてして (29) 式,  $r \in R(f)$ ,  $i \in T$  にてして

$$(32) \quad \tilde{g}_{ct}(\tilde{f}) = (I - \tilde{\delta}(\tilde{f}))^{-1} \tilde{r}(\tilde{f})$$

である。ここで  $\tilde{r}(\tilde{f}) = (\sum_{s \in R(\tilde{f})} \tilde{P}_{rs}(\tilde{f}(r)) \sigma_s^*(s), \sum_{s \in R(\tilde{f})} \tilde{P}_{is}(\tilde{f}_i) \sigma_i^*(s))^T$  である。

(証明)  $Rc(f)UT$  も  $\tilde{\delta}(\tilde{f})$  のもとで閉じた状態集合  $\tilde{T}$  を含むと仮定する。  $\tilde{T}$  に属する状態のうちで最小のレベルをもつものをとれば、必ず  $R(\tilde{f})$  へ遷移する正の確率をもち、矛盾である。

(32) 式は (30) 式から明らかである。

### [定理 2]

(i) 問題 (26) は一意解  $\tilde{g}^* = \tilde{g}(\tilde{f}^*) = \max \{ \tilde{g}(\tilde{f}) \mid \tilde{f} \in \tilde{F} \}$  をもつ。

(ii)  $\{ \bar{\sigma}(r) : r \in R \}$  が  $\bar{\sigma}(r) \geq \sigma^*(r) - \varepsilon_2$  ( $r \in R$ ) を満たすものとす

る, このとき

$$(33) \quad \tilde{g}^*(r) = \max \{ \bar{\sigma}(r), T(r) \bar{g}^* \} \quad (r \in R)$$

$$\bar{g}_i^* = T(i) \bar{g}^* \quad (i \in T)$$

の一意解  $\bar{g}^* = \tilde{g}(\hat{f}^*)$  は  $\tilde{g}^* - \bar{g}^* \leq \varepsilon_2 e$  を満たす。

(証明) (i)  $\hat{f}^*$  は (26) 式右辺を最大化する政策であり,

$$\tilde{g}_s^* = \sigma_s^* \quad (r \in R_s(\hat{f}^*)), \quad \tilde{g}_{ct}^* = T_{ct}(\hat{f}^*) \tilde{g}^* \quad (r \in R_c(\hat{f}^*), i \in T)$$

を満たす。 (29), (30) より  $\tilde{g}^* = \tilde{g}(\hat{f}^*) \leq \max_{\hat{f} \in \hat{F}} \tilde{g}(\hat{f})$  である。 一方,

$$\max_{\hat{f} \in \hat{F}} \tilde{g}(\hat{f}) = \tilde{g}(\hat{f}') \text{ とおけば" (26) より"}$$

$$(34) \quad \tilde{g}_s^* \geq \sigma_s^* = \tilde{g}_s(\hat{f}'), \quad \tilde{g}_{ct}^* \geq T_{ct}(\hat{f}') \tilde{g}^* \quad (r \in R_c(\hat{f}'), i \in T)$$

が " (26)" に満たす。 (34) より  $r \in R_c(\hat{f}')$ ,  $i \in T$  に満たす  $i \in T$

$$(35) \quad \tilde{g}_{ct}^* \geq \tilde{r}(\hat{f}') + \tilde{S}(\hat{f}') \tilde{g}_{ct}^*$$

が " (26)" に満たす。補題 1 より  $\tilde{g}^* \geq \tilde{g}(\hat{f}')$  が成り立つ。

(ii) 停止利得  $\{\bar{g}(r); r \in R\}$  をもつ  $\hat{f}$  の利得を  $\bar{g}(\hat{f})$  で表わせば,

(29) より

$$\bar{g}_s(\hat{f}) = \bar{\sigma}_s = (\bar{\sigma}(r)) \quad (r \in R_s(\hat{f}))$$

である,  $r \in R_c(\hat{f})$ ,  $i \in T$  に満たす  $i \in T$  は (32) より

$$\bar{g}_{ct}(\hat{f}) = (I - \tilde{S}(\hat{f}))^{-1} \bar{r}(\hat{f})$$

である。すなはち  $\bar{r}(\hat{f}) = (\sum_{s \in R_s(\hat{f})} \tilde{P}_{rs}(\hat{f}(r)) \bar{\sigma}(s), \sum_{s \in R_c(\hat{f})} \tilde{P}_{cs}(\hat{f}_c) \bar{\sigma}(s))^T$  である。

$r \in R_s(\hat{f}^*)$  に満たす  $i \in T$

$$\bar{g}_s^* - \tilde{g}_s^* \geq \bar{\sigma}_s - \sigma_s^* \geq -\varepsilon_2 e$$

である,  $r \in R_c(\hat{f}^*)$ ,  $i \in T$  に満たす  $i \in T$  は補題 1 より

$$\begin{aligned} \bar{g}_{ct}^* - \tilde{g}_{ct}^* &\geq \bar{g}_{ct}(\hat{f}^*) - \tilde{g}_{ct}(\hat{f}^*) = (I - \tilde{S}(\hat{f}^*))^{-1} (\bar{r}(\hat{f}^*) - \tilde{r}(\hat{f}^*)) \\ &\geq (I - \tilde{S}(\hat{f}^*))^{-1} \{-\varepsilon_2 (I - \tilde{S}(\hat{f}^*)) e\} = -\varepsilon_2 e \end{aligned}$$

が成り立つ。

## 4. 逐次近似法

部分マルコフ決定過程は、Platzman (1977) 等により提案された逐次近似法を用ひれば、任意の  $\varepsilon_2 > 0$  に対して  $\bar{\sigma}(r) \geq \sigma^*(r) - \varepsilon_2$  となる  $\varepsilon_2$ -最適政策  $\bar{f}(r)$  を有限回の反復で求めることができる。このようにしてえられた  $\{\bar{\sigma}(r); r \in R\}$  を用いて (33) 式を解く逐次近似法として、 $n=0, 1, \dots$  にて

$$(36) \quad \bar{g}^{n+1}(r) = \max \{\bar{\sigma}(r), U(r)\bar{g}^n\} \quad (r \in R)$$

$$\bar{g}_i^{n+1} = U(i)\bar{g}^n_i \quad (i \in T)$$

を考える。ただし初期値として  $\bar{g}^0(r) = \bar{\sigma}(r) \quad (r \in R)$ ,  $\bar{g}_i^0 = \min_r \bar{\sigma}(r)$  ( $i \in T$ ) をとることにする。 $(36)$  式右辺の最大を与える政策を  $\bar{f}^{n+1}$  とする、 $R_s^{n+1} = R_s(\bar{f}^{n+1})$ ,  $R_c^{n+1} = R_c(\bar{f}^{n+1})$  とする。 $r \in R_s^{n+1} \subset T$  のときは  $\bar{g}_s^{n+1} = \bar{\sigma}_s$  である。

## [定理3]

逐次近似法 (36) は (33) 式の解へ単調に収束する。すなわち、 $n=0, 1, 2, \dots$  にて

$$(37) \quad \bar{g}^n \leq \bar{g}^{n+1}, \quad \lim_{n \rightarrow \infty} \bar{g}^n = \bar{g}^*$$

$$(38) \quad R_s^n > R_s^{n+1}, \quad \lim_{n \rightarrow \infty} R_s^n = R_s(\bar{f}^*)$$

がなり  $T$  。

(証明)  $\bar{\sigma}_{\max} = \max_{r \in R} \bar{\sigma}(r)$  とあれば、 $\bar{\sigma}_{\max} e \geq \bar{g}^1 \geq \bar{g}^0$  である。

今  $\bar{\sigma}_{\max} e \geq \bar{g}^n \geq \bar{g}^{n-1}$  であり  $T$  ものとすれば、

$$(39) \quad \bar{\sigma}_{\max} \geq U(r)\bar{g}^n \geq U(r)\bar{g}^{n-1} \quad (r \in R_c^{n+1})$$

$$\bar{g}_{\max} \geq U(i)\bar{g}^n \geq U(i)\bar{g}^{n+1} \quad (i \in T)$$

である。ゆえに (36) 式より

$$\bar{g}_{\max} e \geq \bar{g}^{n+1} \geq \bar{g}^n$$

である。 $\bar{g}^n$  は単調に収束する。 $\bar{g}^n \rightarrow \bar{g}$  とすれば  $\bar{g}$  は (33) 式を満たす。 $\bar{g} = \bar{g}^*$  である。(38) 式及 (39) 式より明らかである。

#### [定理 4]

(i)  $r \in R_c^{n+1}$ ,  $i \in T$  にて  $\bar{g}_{ct}^{n+1} - \bar{g}_{ct}^n = \alpha (\geq 0)$  とおけば。

$$(40) \quad \bar{g}_{ct}(\bar{f}^{n+1}) \geq \bar{g}_{ct}^n + \gamma e = \bar{g}_{ct}^n$$

である。この  $\gamma = \alpha_{\min} / \{1 - \min_{r,i} [\hat{S}(\bar{f}^{n+1})e]_{r,i}\}$  である。

(ii)  $R_s^{n+1} = R_s(\bar{f}^*)$  なる  $T$  のものとする。 $r \in R_c^{n+1}$ ,  $i \in T$  にて

$$U_{ct} \bar{g} - \bar{g}_{ct} \leq b \quad (b > 0 \text{ である})$$

$$(41) \quad \bar{g}_{ct}^* - \bar{g}_{ct} \leq (I - \hat{S}(\bar{f}^*))^{-1} b$$

である。この  $\bar{g}_s = \bar{g}_s$  である、 $\bar{g}_{ct}$  は (40) 式右辺で定義される。

(iii)  $r \in R_s^{n+1}$  にて  $U_s \bar{g} < \bar{g}_s - \varepsilon_1 e$ ,  $r \in R_c^{n+1}$ ,  $i \in T$  にて

$$\bar{g}_{ct}^* - \bar{g}_{ct} \leq \varepsilon_1 e \quad (T \neq \emptyset \text{ ならば } R_s^{n+1} = R_s(\bar{f}^*) \text{ である})$$

(証明) (i)  $r \in R_c^{n+1}$ ,  $i \in T$  にて

$$\begin{aligned} \bar{g}_{ct}(\bar{f}^{n+1}) - \bar{g}_{ct}^n &= (\bar{r}(\bar{f}^{n+1}) + \hat{S}(\bar{f}^{n+1})\bar{g}_{ct}(\bar{f}^{n+1})) - (\bar{r}(\bar{f}^{n+1}) + \hat{S}(\bar{f}^{n+1})\bar{g}_{ct}^n) + \alpha \\ &= \hat{S}(\bar{f}^{n+1})(\bar{g}_{ct}(\bar{f}^{n+1}) - \bar{g}_{ct}^n) + \alpha \end{aligned}$$

であるから、

$$\tilde{g}_{ct}(\bar{f}^{n+1}) - \bar{g}_{ct}^n = (I - \tilde{S}(\bar{f}^{n+1}))^{-1}a = y$$

であります。

$$y = a + \tilde{S}(\bar{f}^{n+1})y \geq (a_{\min} + y_{\min} \cdot \min_{r,i} \{ [\tilde{S}(\bar{f}^{n+1})e]_{r,i} \})e$$

すな (40) 式をうる。

(ii) (33) 式と系 2 から導かれます。

(iii)  $r \in R_s^{n+1} \cap R_c(\bar{f}^*)$  の有無によるものとすれば

$$(42) \quad \bar{\sigma}(r) = \bar{g}(r) \leq \bar{g}^*(r) = U(r)\bar{g}^*$$

すな  $\bar{T}$  は  $T$  の後尾です。後尾すな

$$\begin{aligned} U(r)\bar{g}^* &= \sum_{s \in R_s(\bar{f}^*)} \tilde{P}_{rs}(\bar{f}^*)\bar{\sigma}(s) + \sum_{s \in R_s^{n+1} \cap R_c(\bar{f}^*)} \tilde{P}_{rs}(\bar{f}^*)\bar{g}^*(s) + \sum_{s \in R_c^{n+1}} \tilde{P}_{rs}(\bar{f}^*)\bar{g}^*(s) \\ &\quad + \sum_{j \in T} \tilde{P}_{rj}(\bar{f}^*)\bar{g}_j^* \\ &\leq \sum_{s \in R_s} \tilde{P}_{rs}\bar{\sigma}(s) + \sum_{s \in R_s^{n+1} \cap R_c} \tilde{P}_{rs}\bar{g}^*(s) + \sum_{s \in R_c^{n+1}} \tilde{P}_{rs}(\bar{g}(s) + \varepsilon_2) + \sum_{j \in T} \tilde{P}_{rj}(\bar{g}_j + \varepsilon_1) \\ &< \bar{\sigma}(r) + \sum_{s \in R_s^{n+1} \cap R_c} \tilde{P}_{rs}(\bar{g}^*(s) - \bar{\sigma}(s)) - \varepsilon_1 \sum_{s \in R_s^{n+1}} \tilde{P}_{rs} \end{aligned}$$

であります、 $c = \max_{r \in R_s^{n+1} \cap R_c} \bar{g}^*(r) - \bar{\sigma}(r)$  とおけば  $c(1 - \sum_{s \in R_s^{n+1} \cap R_c} \tilde{P}_{rs}) < -\varepsilon_1 \sum_{s \in R_s^{n+1}} \tilde{P}_{rs}(\bar{f}^*)$

となり矛盾である。

(41) 式の評価には  $(I - \tilde{S}(\bar{f}^*))^{-1}b$  の上限 (下限) が必要であります。

$$(I - \tilde{S}(\bar{f}^*))^{-1}b = \sum_{m=0}^{\infty} \tilde{S}(\bar{f}^*)^m b \leq \max_{f \in F} \sum_{m=0}^{\infty} \tilde{S}(f)^m b$$

でありますから、上限  $y$  は次のマルコフ決定過程に帰着されます。

$$(43) \quad y = \max_{f \in F} \{ b + \tilde{S}(f)y \}.$$

この問題を解く逐次近似法として  $n = 0, 1, \dots, l = T, \dots, 0$

$$(44) \quad y^{n+1} = \max_f \{ b + \tilde{S}(f)y^n \}, \quad y^0 = b$$

$$(45) \quad \bar{t}^{n+1} = \max_f \bar{s}(f) \bar{t}^n, \quad \bar{t}^0 = e$$

を考へる。補題 1 より  $N \leq |R_e^n| + |T|$  を満たすある  $N$  で

$$(46) \quad \bar{t} = \max_{r,i} t_{r,i}^N < 1$$

とする、 $\bar{y} = \max_{r,i} y_{r,i}^{N-1}$  とかければ

$$(47) \quad \bar{y} \leq \frac{\bar{y}}{1-\bar{t}} e$$

である。

以上の結果から多電連鎖マルコフ決定過程  $(I, r, p, f)$  の  $\varepsilon$  最適政策をもとめる次の逐次近似法がえられる。

1. アルゴリズム (17) に より  $C(r)$  ( $r \in R$ ),  $D_i$  ( $i \in C(r)$ ),  $T$  をもとめる。

2. 各  $r \in R$  にたいする部分マルコフ決定過程  $(C(r), r, p, f(r))$  を逐次近似法 ("たとえば" Platzman (1977)) に より解く。判得  $\bar{\sigma}(r)$ ,  $\varepsilon_2$  最適政策  $\bar{f}(r)$  を求める。

3.  $\bar{g}^0(r) = \bar{\sigma}(r)$  ( $r \in R$ ),  $\bar{g}_i^0 = \min_r \bar{\sigma}(r)$  ( $i \in T$ ),  $n=0$  とおく。

4. (36) 式から  $\bar{g}^{n+1}$  を計算し,  $r \in R_e^{n+1}$ ,  $i \in T$  に  $T = \cup_i T_i$

$$\bar{g}_{ct}^{n+1} - \bar{g}_{ct}^n = \alpha \leq \delta e$$

となるれば 5. へ。さもなくば  $m=n+1$  として (36) 式の計算にもどる。

5. (40) 式右辺の  $\bar{g}_{ct}$  を計算し,  $\varepsilon_1 = \varepsilon - \varepsilon_2 (> 0)$  に  $T = \cup_i T_i$

$$U(r) \bar{g} < \bar{\sigma}(r) + \varepsilon_1 \quad (r \in R_e^{n+1})$$

となるれば 6. へ。さもなくば  $n=n+1$ ,  $\delta = \alpha \delta$  ( $0 < \alpha < 1$ )

とし て 4. へ。

6. (44) ~ (47) 式より  $\bar{t}, \bar{y}$  をもとめ,

$$\frac{\bar{y}}{1-\bar{t}} \leq \varepsilon_1$$

とすれば 7. へ。でもなければ  $n=n+1, \delta=\alpha\delta$  とし て 4.  
へ。

7.  $r \in R_s^{n+1}$  に て は  $\bar{f}_i = \bar{f}_i(r)$  ( $i \in C(r)$ ),  $i \in T = T \cup i$   
に て  $\bar{f}_i = \bar{F}_i^{n+1}$  とあく。各  $r \in R_c^{n+1}$  に て,  $\bar{f}^{n+1}(r) =$   
( $i^*, k^*$ ) とすば  $\bar{f}_{i^*} = k^*$  とすき,  $i (\neq i^*) \in C(r)$  に て  $\bar{f}_i$   
は  $i^*$  を吸収状態とする政策  $\bar{f}(r) \in F(r)$  を求め,  $\bar{f}_i = \bar{f}_i(r)$   
とあく。政策  $\bar{f}$  は  $\varepsilon$ -最適政策である。

### 参考文献

- Bather, J. (1973), "Optimal Decision Procedures for Finite Markov chains, III," Adv. Appl. Prob. 5, 541 - 553.
- Federgrun, A. and P. J. Schweitzer (1984), "A Fixed Point Approach to Undiscounted Markov Renewal Programs," SIAM J. Alg. Dis. Math. 5, 539 - 550.
- Platzman, L. (1977), "Improved Conditions for Convergence in Undiscounted Markov Renewal Programs," Opns. Res. 25, 529 - 533.
- Schweitzer, P. J. (1984), "A Value-Iteration Scheme for Undiscounted Multi-chain Markov Renewal Programs," Zeit. Opns. Res. 28, 143 - 152.
- Whittle, P. (1983), Optimization over Time, Vol. 2, John Wiley, Chichester.