

時間平均多重連鎖セミ・マルコフ決定過程における 修正政策及復法の収束について

甲南大学 理学部 大野勝久 (Katsuhisa Ohno)

1. はじめに

有限状態、有限決定セミ・マルコフ決定過程にお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、政策及復法 (PIM)、逐次近似法 (SAM)、線形計画法 (LP)、修正政策及復法 (MPIM) が知られてる [1]。特に PIM はよく知られてるが、PIM, LP は共に多状態問題にたっては適用が困難であり、これら問題にたっては SAM およびその一般化としての MPIM が有力な手法として研究されてきた。しかし SAM, MPIM は多重連鎖問題にたっては無力であった。最近 Schweitzer [2] は多重連鎖問題にたって SAM を提案し、その収束を示してゐる。しかしながらそのアルゴリズムは複雑であり、実用的とは思われない。本論文ではより一般的に MPIM の収束について論ずる。

2. 政策反復法

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 でとりうる決定の有限集合

$r_{ij}(k) (i \in I, k \in K_i)$: i において k をとったとき j へ平均利得

$P_{ij}(k) (i, j \in I, k \in K_i)$: i において k をとったとき j へ遷移する確率

$T_i(k) (i \in I, k \in K_i)$: 平均遷移時間 ($T_i(k) > 0$ とする)

F : 定常政策 $f = (f_1, f_2, \dots, f_M)$ の集合 ($f_i \in K_i, i \in I$)

$r(f) = (r_1(f_1), r_2(f_2), \dots, r_M(f_M))^T$ (T は転置を表す)

$P(f) = (P_{ij}(f_i))$: f による遷移確率行列

$T(f) = \text{diag}(T_i(f_i))$

$g(f)$: 定常政策 f のゲイン

$v(f)$: 定常政策 f の相対値

$g(f), v(f)$ は, $P(f)$ の各エルゴード部分連鎖に属する 1 つの状態 s について $v_s = 0$ とおいて, 次の連立一次方程式を解いてえられる。

$$g = P(f)g, \quad v = r(f) + P(f)v - T(f)g \quad (1)$$

ここで, $0 < \alpha < 1$,

$$0 < \tau \leq (1-\alpha) \min_i \left\{ T_i(k) / (1-P_{ii}(k)) ; P_{ii}(k) < 1 \right\} \quad (2)$$

をみたす任意の α , τ にたいして

$$\Omega(f) = \tau T(f)^{-1} \quad (3)$$

とかき, Schweitzer の data 変換:

$$\tilde{P}(f) = I - \Omega(f) + \Omega(f)P(f), \quad \tilde{r}(f) = \Omega(f)r(f) \quad (4)$$

を用ひれば" (1) 式は

$$g = \tilde{P}(f)g, \quad v = \tilde{r}(f) + \tilde{P}(f)v - g \quad (5)$$

となる。ただし、(4) 式の I は単位行列である, (5) 式の解 \tilde{g} , \tilde{v} は $\tilde{g} = \tau g(f)$, $\tilde{v} = v(f)$ である, 任意の $f \in F$ にたいして

$$\tilde{P}(f) \geq \alpha I \quad (6)$$

をみたす。すなわち, セミ・マルコフ決定過程は上記 data 変換で任意の $f \in F$ にたいして (6) 式をみたす同値なマルコフ決定過程に変換される。

以下記号を簡略化するため変換されたマルコフ決定過程にたいして再び $r(f)$, $P(f)$ を用ひることにする。すなわち

$$g = P(f)g, \quad v = r(f) + P(f)v - g \quad (7)$$

の解が $g(f)$, $v(f)$ であり, $P(f)$ は任意の $f \in F$ にたいして

$$P(f) \geq \alpha I \quad (8)$$

をみたす。また最大ゲイン $g^* = g(f^*)$, $v^* = v(f^*)$ のみたす最適方程式は

$$g_i^* = \max_{k \in K_i} \left\{ \sum_{j \in I} p_{ij}(k) g_j^* \right\} \quad (i \in I) \quad (9)$$

$$g_i^* + v_i^* = \max_{k \in L_i} \{ r_i(k) + \sum_{j \in I} p_{ij}(k) v_j^* \} \quad (10)$$

で与えられる。ここで $L_i = \{ k \in K_i ; (9) \text{ 式右辺を最大化する } k \} \text{ である}。$

[補題 1]

$P(f^*)$ のエルゴード部分連鎖の状態集合を E_r^* ($r=1, \dots, R^*$)、過渡状態の集合を T^* とし、 $g_r^* = g_i(f^*)$ ($i \in E_r^*$) とおく。

$$g_1^* < g_2^* < \dots < g_R^*$$

ならば、

i) $T^* \ni i$ にたって $g_1^* \leq g_i(f^*) \leq g_R^*$ である。

ii) $E_r^* \ni i$ にたって

$$f_i^* \in D_i = K_i - \{ k \in K_i ; E_1^* \cup \dots \cup E_{r-1}^* \cup \{ i \in T^* ; g_i(f^*) \leq g_r^* \} \ni j \\ \text{にたって } p_{ij}(k) > 0 \}$$

であり、 E_r^* は任意の $f_i^* \in D_i$ ($i \in E_r^*$) で閉じている。

(証明) i) は明らかであり、 $f_i^* \notin D_i$ ($i \in E_r^*$) ならば E_r^* がエルゴード部分連鎖であることを矛盾するとかく ii) の前半が示されると、 E_r^* が閉じてなら $f_i^* \in D_i$ ($i \in E_r^*$) が存在したとすれば、(9) 式より

$$g_r^* = \max_{k \in K_i} \left\{ \sum_{j \in I} p_{ij}(k) g_j(f^*) \right\} \geq g_r^* \sum_{j \in E_r^*} p_{ij}(f_i^*) + \sum_{j \notin E_r^*} p_{ij}(f_i^*) g_j(f^*) > g_r^*$$

となり、矛盾である。

多重連鎖マルコフ決定過程の最適定常政策 f^* を求める最もよく知られた手法は次の Howard の政策反復法である。

1. 初期政策 f^0 を与え, $m=0$ とおく。

2. (値決定ルーテン)

$P(f^n)$ の各エルゴード部分連鎖に属する 1 つの状態 s で

$v_s = 0$ とかき, (7) 式を解いて $g(f^n)$, $v(f^n)$ を定める。

3. (政策改良ルーテン) 各 $i \in I$ で

$$L_i^{n+1} = \{ \sum_{j \in I} P_{ij}(k) g_j(f^n) \text{ を最大化する } k \in K_i \}$$

$$F_i^{n+1} = \{ r_i(k) + \sum_{j \in I} P_{ij}(k) v_j(f^n) \text{ を最大化する } k \in L_i^{n+1} \}$$

を定め, $f_i^{n+1} \in F_i^{n+1}$ ならば $f_i^{n+1} = f_i^n$ とかき, さもなければ f_i^{n+1} を F_i^{n+1} の適当な要素とする。全ての $i \in I$ で $f_i^{n+1} = f_i^n$ となれば停止。 f^{n+1} は最適政策 f^* であり, 最大ゲイン $g^* = g(f^n)$, 相対値 $v^* = v(f^n)$ である。さもなければ $m = m+1$ とおいてステップ 2 へ。

3. 修正政策及復法

PIM の値決定ルーテンと有限回の逐次近似法を組みかえた手法が MPIM である。下記にたてて初期ベクトル w^0 ではじまる次の逐次近似法 $\{ w^l ; l = 0, 1, \dots \}$ を考える。

$$w^{l+1} = r(f) + p(f) w^l. \quad (11)$$

このとき, $l = 0, 1, \dots$ にて

$$w^l = l g(f) + v(f) + p(f)^l (w^0 - v(f)) \quad (12)$$

すり下ち [3], 次の補題が成立する。

[補題 2]

$P(f)$ のエルゴード部分連鎖の状態集合を $E_r (r=1, \dots, R)$, 過渡状態の集合を T とし , $g_r = g_i(f) (i \in E_r)$, $g_1 < g_2 < \dots < g_R$ とする。このとき

$$\text{i) } \lim_{l \rightarrow \infty} w_i^{l+1} - w_i^l = g_i(f) \quad (13)$$

ii) $E_r \ni i, s_r \in E_r \cup T$, $v_{s_r}(f) = 0$ であれば。

$$\lim_{l \rightarrow \infty} w_i^l - w_{s_r}^l = v_i(f) \quad (14)$$

でみる , $s_p \in E_p (p < r)$ であれば。

$$\lim_{l \rightarrow \infty} w_i^l - w_{s_p}^l = \infty . \quad (15)$$

$$\text{iii) } P(f) = \begin{pmatrix} Q(f) & 0 \\ R(f) & S(f) \end{pmatrix} \quad (16)$$

とする。ここで $Q(f)$ は $E_1 \cup \dots \cup E_R$ に対応し , $R(f), S(f)$ は T に対応する。このとき

$$v_i^l = w_i^l - w_{s_r}^l \quad (i \in E_r, r=1, \dots, R) \quad (17)$$

$$g_i^l = w_i^l - w_i^{l-1} \quad (i \in T) \quad (18)$$

とすると , r_t, g_t が T に対応する部分ベクトルを表すと ,

$$u_i^{l+1} (i \in T) \in u_i^1 = w_i^1 \quad (i \in T) \text{ から}$$

$$u_i^{l+1} = r_t(f) + R(f)v_i^l - g_t^l + S(f)u_i^l \quad (19)$$

で求めれば、 $T \ni i$ に対して

$$\lim_{l \rightarrow \infty} u_i^l = v_i(f) . \quad (20)$$

さて T が $T \rightarrow$ 。

(証明) (8) 式より

$$\lim_{l \rightarrow \infty} P(f)^l = P^*(f) = \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{k=1}^N P(f)^l, \quad P^*(f)P(f) = P(f)P^*(f) = P^*(f) \quad (21)$$

であるから、i) は (12) 式より明かである。ii) も

$E_r \ni i \in T_i \cup \{r\}$

$$g_r = \lim_{l \rightarrow \infty} w_i^{l+1} - w_i^l = \lim_{l \rightarrow \infty} v_i(f_i) + \sum_{j \in I} p_{ij}(f_i)(w_j^l - w_{sj}^l) - (w_i^l - w_{sr}^l) \quad (22)$$

である。(7) 式がなりたつことから (14) 式がえられ。

(15) 式は (12) 式から

$$\lim_{l \rightarrow \infty} w_i^l - w_{sp}^l = l(g_r - g_p) + v_i(f) + P(f)^l (w^0 - v(f))_i - P(f)^l (w^0 - v(f))_{sp}$$

となることから導かれる。ii) の証明同様、(17), (18) 式の
 v_i^l, g_i^l は $l \rightarrow \infty$ で $v_i(f)$ ($i \in E_r$), $g_i(f)$ へ収束し、 $S(f)^l \rightarrow 0$ ($l \rightarrow \infty$)
> こと (20) 式がえられる。

PIM の値決定ルーチンに上記逐次近似法を用いれば MPIM
> がえられる。(しかし問題 1 で述べた一般的な多重連鎖問題に
> たいする MPIM は複雑であり、以下では簡単のため $R^* \leq 2$,
> すなわち $P(f^*)$ がただか 2 つのエルゴード部分連鎖をもつ
> ことが知られた問題にたいする MPIM を述べ、次節でその収
> 束を示す。

修正政策反復法 ($R^* \leq 2$)

1. 初期ベクトル v^0 , 非負整数 m , 大きな数 L , 小さな正数
 ε, δ を定め, $m = c = 0$, $E_1 = I$, $T = \emptyset$, $R^* = 1$, $D_i = K_i$ ($i \in I$)
> とおく。

2. (政策改良ルーチン)

i) $i \in E_r$ ($1 \leq r \leq R^*$) にたゞ $i \in T$

$$x_i^{n+1} = \max_{k \in D_i} \left\{ r_i(k) + \sum_{j \in I} p_{ij}(k) v_j^n - v_i^n \right\} \quad (23)$$

を計算し, f_i^n が最大値を与えれば $f_i^{n+1} = f_i^n$ とおく, さもなくば f_i^{n+1} を最大値を与える任意の $k \in D_i$ とする。

$$s_r = \arg \min_{i \in E_r} \{ x_i^{n+1} \} \quad \text{とおく}.$$

ii) $i \in T$ にたゞ $i \in T$

$$L_i^{n+1} = \{ k \in K_i : \sum_{j \in I} p_{ij}(k) g_j^m \geq \max_{k \in K_i} \{ \sum_{j \in I} p_{ij}(k) g_j^m \} - \delta \} \quad (24)$$

$$x_i^{n+1} = \max_{k \in L_i^{n+1}} \left\{ r_i(k) + \sum_{j \in I} p_{ij}(k) v_j^n - v_i^n \right\} \quad (25)$$

を計算し, f_i^n が最大値を与えれば $f_i^{n+1} = f_i^n$ とおく, さもなくば f_i^{n+1} を最大値を与える任意の $k \in L_i^{n+1}$ とする。

3. (値決定ルーティン)

$$w^0 = v^n + x^{n+1} \quad (26)$$

とおく, $l = 0, 1, \dots, m-1$ にたゞ $i \in T$

$$w^{l+1} = r(f^{n+1}) + p(f^{n+1}) w^l \quad (27)$$

を計算し,

$$v_i^{n+1} = w_i^m - w_{s_r}^m \quad (i \in E_r, 1 \leq r \leq R^*) \quad (28)$$

$$g_i^m = w_i^m - w_i^{m-1} \quad (i \in I) \quad (29)$$

とおく。 $u_i^l = w_i^l$ ($i \in T$) とおく, $l = 0, \dots, m-1$ にたゞ $i \in T$

$$u^{l+1} = r_t(f^{n+1}) + R(f^{n+1}) v^{n+1} - g_t^m + S(f^{n+1}) u^l \quad (30)$$

を計算し, v_i^{n+1} ($i \in T$) を次式で定める。

$$v_i^{n+1} = u_i^m \quad (31)$$

4. $m = n+1$ をみる, $c = 1$ ならばステップ° 2へ. $c = 0$ ならば
は " $\|v^n\|_d \equiv \max_{i \in I} v_i^n - \min_{i \in I} v_i^n$ を求め, $\|v^n\|_d < L$ ならば
ステップ° 2へ. さもなければステップ° 5へ.

5. (状態分類ルーチン)

s_1 を含む全ての $f_i \in K_i$ で閉じてなる状態の集合 E'_1 および
 $\arg \max_{i \in I} v_i^n$ を含む f^n で閉じてなる状態の集合 E'_2 を定める。

$$\max_{i \in E'_1} x_i^n - \min_{i \in E'_1} x_i^n < \varepsilon \quad \text{とし, } R^* = 2, \quad E_r = E'_r \quad (r=1,2),$$

$$T = I - E_1 - E_2, \quad D_i = K_i - \{k \in K_i; j \notin E_2 \text{ にたどり } R_j(k) > 0\}$$

($i \in E_2$), $c = 1$ をみる, ステップ° 2へ. さもなければ、

$L = 2L$ をみてステップ° 2へ.

4. 修正政策及復法の収束

まず $c = 0$, すなわち $E_1 = I$, $T = \emptyset$, $D_i = K_i \quad (i \in I)$ にたどり
ある収束を論ずる。 (7) 式より, (23), (28) 式は各々

$$x^{n+1} = g(f^{n+1}) - (I - P(f^{n+1}))(v^n - v(f^{n+1})) \quad (32)$$

$$v^{n+1} = (m+1)g(f^{n+1}) + v(f^{n+1}) + P(f^{n+1})^{n+1}(v^n - v(f^{n+1})) - w_s^m e \quad (33)$$

と書くことができる。ここで $e = (1, \dots, 1)^T$ である。

[補題 3]

$$i) \quad x^{n+1} = d^{n+1} + P(f^n)^{n+1} x^n \quad (34)$$

$$\begin{aligned} d^{n+1} &= \max_{f \in F} \{(m+1)[P(f)g(f^n) - g(f^n)] + r(f) + P(f)v(f^n) - g(f^n) - v(f^n) \\ &\quad + (P(f) - P(f^n))P(f^n)^{n+1}(v^{n+1} - v(f^n))\} \geq 0 \end{aligned} \quad (35)$$

$d_i^{n+1} = 0$ となる必要十分条件は $f_i^{n+1} = f_i^n$ である。

$$\text{ii) } x^{n+1} \geq p(f^n)^{m+1} x^n, \quad x^{n+1} + v^n \geq r(f^*) + p(f^*) v^n \quad (36)$$

$$p^*(f^*) x^{n+1} \geq g(f^*) \geq g(f^{n+1}) = p^*(f^{n+1}) x^{n+1}, \quad p^*(f^n) x^{n+1} \geq g(f^n) \quad (37)$$

$$\text{iii) } \Delta_n = \max_{i \in I} x_i^n, \quad \nabla_n = \min_{i \in I} x_i^n \quad (38)$$

とかいてば

$$\Delta_n e \geq g(f^*) \geq g(f^{n+1}) \geq \nabla_n e \geq \nabla_n e \quad (39)$$

iv) 状態の集合 C が $p(f)^{m+1}$ で閉じてがあれば、 $p(f)$ でも閉じていい。

(証明) i) (12) 式と同様

$$v^n = mg(f^n) + v(f^n) + p(f^n)^m (v^{n+1} + x^n - v(f^n)) - w_s^m e$$

さて、この式と (33) 式を (23) 式に代入し、(7), (32) 式を用いて整理すれば i) が示される。ii) は i), (21) 式より導かれ、iii), iv) は明らかである。

[補題 4]

i) $\bar{v} = \lim_{n \rightarrow \infty} v_n$ が存在し、空でない集合 $C \subset I$ に対して $\lim_{n \rightarrow \infty} x_i^n = \bar{v} \quad (i \in C)$

ii) 有限な N が存在し、全ての $n \geq N$ に対して

$$f_i^n = f_i \in K_i, \quad d_i^n = 0 \quad (i \in C) \quad (41)$$

であり、 C は $p(f^n)$ で閉じてある。

(証明) (39) 式より v_n は \bar{v} へ収束する。 $i \in I$ に対して、 $y_i = \liminf_{n \rightarrow \infty} x_i^n$, $z_i = \limsup_{n \rightarrow \infty} x_i^n$ とかく、
 $A = \{i \in I ; \text{無限に多くの } n \text{ で } x_i^n = v_n\}$

とおけば A は空でなく, $A \ni i$ にたいして $y_i = \bar{V}$ が成り立つ。

したがって

$$C = \{ i \in I ; y_i = \bar{V} \}$$

とおけば C は空でない。定義より $i \in C$ にたいして $\{n\}$ の部

分列 $L(i) = \{l_i\}$, $U(i) = \{u_i\}$ が存在す。

$$\lim_{l \rightarrow \infty} x_i^l = \bar{V}, \quad \lim_{u \rightarrow \infty} x_i^u = z_i$$

である。部分列 $\{n_k; k = 0, 1, \dots\}$ を

$$n_{2k} \in U(i), \quad n_{2k+1} \in L(i), \quad n_{2k+1} < n_{2k} < n_{2k+1}, \quad n_{2k+1} - n_{2k} < \infty$$

構成し, $P(2k) = p(f^{n_{2k+1}-1})^{m+1} \dots p(f^{n_{2k}})^{m+1}$ とおけば, (36)

式より

$$x^{n_{2k+1}} \geq p(f^{n_{2k+1}-1})^{m+1} x^{n_{2k+1}-1} \geq \dots \geq P(2k) x^{n_{2k}}$$

である。ゆえに

$$x_i^{n_{2k+1}} \geq p_i(2k) x_i^{n_{2k}} + \sum_{j \neq i} p_{ij}(2k) x_j^{n_{2k}} \geq V_{n_{2k}} + p_i(2k)(x_i^{n_{2k}} - V_{n_{2k}})$$

である。したがって, (8) 式より $p_i(2k) > 0$ であるから $z_i \leq \bar{V}$ となる。

ii) $C \neq i$ にたいして十分大きい n で

$$x_i^n > \bar{V} + \gamma \tag{42}$$

とする $\gamma > 0$ が存在する。 (36) 式より $i \in C$ にたいして

$$x_j^{n+1} \geq V_n + \sum_{j \notin C} p_{ij}(f^n)^{m+1} (x_j^n - V_n) \geq V_n + \gamma \sum_{j \notin C} p_{ij}(f^n)^{m+1}$$

である。ゆえに i) が $\sum_{j \notin C} p_{ij}(f^n)^{m+1} = 0$ である, 補題 3-iv)

より C は $p(f^n)$ で閉じてなる。 (\Rightarrow が成り立つ) (34) 式より

$$\lim_{n \rightarrow \infty} d_i^{n+1} = \lim_{n \rightarrow \infty} (x_i^{n+1} - \sum_{j \in C} p_{ij}(f^n)^{n+1} x_j^n) = 0$$

と \bar{f}), 補題 3-i) より f_i^n は必ず決良 $\bar{f}_i \in K_i$ へ収束する。C, K_i もともに有限であるから ii) もなりたつ。

[定理 1]

$C = I$ ならば、 $I \ni x_i = \bar{x}_i$ にて

$$f_i^n = f_i^* (n \geq N), \quad \lim_{n \rightarrow \infty} x_i^n = g_i^*, \quad \lim_{n \rightarrow \infty} v_i^n = v_i(f^*)$$

がなりたつ。

(証明) 補題 3-ii) より $v_i(\bar{f}) = v_i(f^*)$

$$g_i(f^*) = g_i^* e = g_i(\bar{f}) = \bar{v}_i e \quad (43)$$

がえられる。任意の $\varepsilon > 0$ にて $\bar{N} (\geq N)$ が存在し、
 $n \geq \bar{N}$ において

$$x^n = \bar{v}_i e + \varepsilon_i^n, \quad |\varepsilon_i^n| < \varepsilon \quad (44)$$

であるから

$$g_i(\bar{f}) + v_i^n + \varepsilon_i^n = r(\bar{f}) + p(\bar{f}) v_i^n$$

である、(7) 式よ'

$$\lim_{n \rightarrow \infty} v_i^n = v_i(\bar{f}) \quad (45)$$

がなりたつ。ゆえに $d^n = 0$ よ'

$$g_i(\bar{f}) + v_i(\bar{f}) = \max_{f \in F_i} \{ r(f) + p(f) v_i(f) \}$$

と $\bar{f} = f^*$ が示される。

以下 $C \neq I$ の場合を考える。 $T = I - C$ とおけば補題 4 より $m \geq N$ にて $p(f^n)$ は

$$P(f^n) = \begin{pmatrix} Q(\bar{f}) & 0 \\ R(f^n) & S(f^n) \end{pmatrix} \quad (46)$$

と表わすことができます。ここで R, S は T に対応する行列です
みる。以下添字 c, t で C, T に対応する部分ベクトルを表す。
すなはち

[定理 2]

$C \neq I$ ならば $R^* = 2$, $g_1^* < g_2^*$, $C = E_1^*$ であり, $E_1^* \ni i$ にて
して $f_i^n = f_i^*$ ($n \geq N$), $\lim_{n \rightarrow \infty} x_i^n = g_i^*$, $\lim_{n \rightarrow \infty} v_i^n = v_i(f^*)$ が
り $I = \bar{f}$ 。

(証明) 前題 3, 4 より

$$g_C(f^n) = Q^*(\bar{f})x_C^n = Q^*(\bar{f})v_C(\bar{f}) = g_C(\bar{f}) = \bar{v}e_C \quad (47)$$

が $I \neq \bar{f}$, $m \geq \bar{N}$ で

$$x_C^n = \bar{v}e_C + \varepsilon_C^n, \quad |\varepsilon_C^n| < \varepsilon \quad (i \in C) \quad (48)$$

である。ゆえに定理 1 の証明と同様にして (45) 式が \neq で
す,

$$v_C^n = v_C(\bar{f}) + \delta_C^n = v_C(f^{n+1}) + \delta_C^n, \quad |\delta_C^n| < \varepsilon \quad (i \in C) \quad (49)$$

である。(32) 式から

$$x_t^{n+1} = g_t(f^{n+1}) - R(f^{n+1})(v_C^n - v_C(f^{n+1})) - (I_t - S(f^{n+1}))(v_t^n - v_t(f^{n+1}))$$

であるから (49) 式より

$$x_t^{n+1} - g_t(f^{n+1}) + R(f^{n+1})\delta_C^n = (I_t - S(f^{n+1}))(v_t(f^{n+1}) - v_t^n) \quad (50)$$

と $\bar{T} \neq I$ で

$$(I_t - S(f^m)^{m+1})(v_t(f^{m+1}) - v_t^m) = \sum_{\ell=0}^m S(f^{m+1})^\ell \{ x_t^{m+1} - g_t(f^{m+1}) + R(f^{m+1}) \delta_c^n \} \quad (51)$$

ゆえに $T = \emptyset$ 。一方、(33) 式は (46), (47), (49) 式を代入すれば

$$v_c^{m+1} = (m+1) \bar{V} e_c + v_c(\bar{f}) + Q(\bar{f})^{m+1} \delta_c^n - w_{S_1}^m e_c \quad (52)$$

$$\begin{aligned} v_t^{m+1} &= (m+1) g_t(f^{m+1}) + v_t(f^{m+1}) + R(f^{m+1})^{m+1} \delta_c^n \\ &\quad + S(f^{m+1})^{m+1} (v_t^m - v_t(f^{m+1})) - w_{S_1}^m e_t \end{aligned} \quad (53)$$

と $T \neq \emptyset$ 。ここで $R(f)^{m+1} = \sum_{\ell=0}^m S(f)^\ell R(f) Q(\bar{f})^{m-\ell}$ である。 $v_{S_1}^{m+1} = v_{S_1}(\bar{f})$
 $= 0$ であるから (52) 式より $w_{S_1}^m = (m+1) \bar{V} + (Q(\bar{f})^{m+1} \delta_c^n)_{S_1}$ である
 し、(53) 式は

$$\begin{aligned} v_t^{m+1} - v_t(f^{m+1}) &= (m+1)(g_t(f^{m+1}) - \bar{V} e_t) + R(f^{m+1})^{m+1} \delta_c^n - (Q(\bar{f})^{m+1} \delta_c^n)_{S_1} e_t \\ &\quad + S(f^{m+1})^{m+1} (v_t^m - v_t(f^{m+1})) \end{aligned}$$

と $T \neq \emptyset$ 。したがって、(42), (51) 式を用いて整理すれば

$$\begin{aligned} v_t^{m+1} - v_t^m &= \sum_{\ell=0}^m S(f^{m+1})^\ell \{ x_t^{m+1} - \bar{V} e_t + R(f^{m+1})(I_c + Q(\bar{f})^{m-\ell}) \delta_c^n \} \\ &\quad - (Q(\bar{f})^{m+1} \delta_c^n)_{S_1} e_t \\ &> (\gamma - 3\varepsilon) e_t + \gamma \sum_{\ell=1}^m S(f^{m+1})^\ell e_t + 2\varepsilon S(f^{m+1})^{m+1} e_t \end{aligned} \quad (54)$$

ゆえに $T \neq \emptyset$ で、 $\gamma > 3\varepsilon$ とすれば

$$\lim_{n \rightarrow \infty} v_i^n = \infty \quad (i \in T) \quad (55)$$

ゆえに $\exists i \in C$ で $R_j(f_i^*) > 0$ ($j \in T$) と仮定すれば、
 (36) 式の右辺が ∞ となる矛盾である。すなわち左辺の $i \in T$ で C は閉じてある。

$$g_c(\bar{f}) = \max_{f \in F} \{ P_c(f) g_c(f) \}, \quad g_c(\bar{f}) + v_c(\bar{f}) = \max_{f \in F} \{ R_c(f) + P_c(f) v_c(f) \}$$

がなりたつ。ゆえに $R^* = 2$, $g_1^* < g_2^*$, $C = E_1^*$, $\bar{f}_i = f_i^* (i \in E_1^*)$ が示される。

定理 1, 2 より $R^* \leq 2$ にたいする MPIM の収束をみる条件のもとで示すこととする。

[定理 3]

修正政策及復法は, $R^* = 1$ または $R^* = 2$, $g_1^* = g_2^*$ ならば常に, $R^* = 2$, $g_1^* < g_2^*$ ならば

i) \bar{N} の存在し, E_2^* は $n \geq \bar{N}$ となる全ての n において $P(f^n)$ で閉じてなる。

ii) $T^* \ni n \mapsto i \in \mathbb{Z}$, $\sum_{j \in T^*} P_{ij}(f_i^n) < 1$ ならば

の条件のもとで、任意の v^0 , $m \in \mathbb{N}$ にて収束する。

(証明) 定理 1, 2 より $R^* = 1$ または $R^* = 2$, $g_1^* = g_2^*$ ならば $C = I$ となり, MPIM は任意の v^0 , m で収束する。 $R^* = 2$, $g_1^* < g_2^*$ ならば定理 2 より $C = E_1^*$ である, 仮定 i) から $m \geq \bar{N}$ において $f_i^n = f_i^* (i \in E_1^*)$, $\sum_{j \in E_2^*} P_{ij}(f_i^n) = 1 (i \in E_2^*)$ がなりたつ。さし $i \in E_1^* \cap \mathbb{Z}$ にて

$$\lim_{n \rightarrow \infty} x_i^n = g_1^*, \quad \lim_{n \rightarrow \infty} v_i^n = v_i(f^*)$$

がなりたつ。仮定 i) より $m \geq \bar{N}$ において $x_j^n (j \in E_2^*)$ は $K_2 \cup D_2 = K_2 - \{k \in K_2 ; j \notin E_2^* \text{ にて } P_{kj}(f^*) > 0\}$ に限定され、それともへと一致する。したがって定理 1 で, $I = E_2^*$ とあれば、

$$\lim_{n \rightarrow \infty} x_i^n = g_2^* \quad (i \in E_2^*)$$

さて、 $T \rightarrow \infty$ のとき $\|v^n\|_d \rightarrow \infty$, $\max_{i \in E_2^*} x_i^n - \min_{i \in E_2^*} x_i^n \rightarrow 0$ ($n \rightarrow \infty$) であるから、ある有限な $m = N'$ でステップの条件が成立し、 $E_1 = E_1^*$, E_2 が定められる。さて、仮定 i) より $E_2 \subset E_2^*$ である、(37) 式より $i \in E_2$ にたいして

$$\sum_{j \in E_2} p_{ij}^*(t) \{(\max_{i \in E_2} x_i^n) - x_j^n\} \geq 0$$

であるから、 $\sum_{j \in E_2} p_{ij}^*(t) = 0$ となる $E_2 > E_2^*$ がなりたつ。ゆえに $E_2 = E_2^*$ であり、 $i \in E_2^*$ にたいして収束性定理 1 より導かれ、 $i \in T = T^*$ にたいして収束は、仮定 ii) より K_i を $\sum_{j \in T^*} p_{ij}(t) < 1$ とするために限らざる。割引利得問題にたいして収束証明 [4] と同様にして示すことができる。

参考文献

[1] 大野, マルコフ決定過程の計算アルゴリズム. 第4回数理計画シンポジウム論文集 pp. 143-158, 1983.

2 Schweitzer, P.J., "A value-iteration scheme for undiscounted multichain Markov renewal programs," Zeitschrift für Operations Res. 28, 143-152, 1984.

3 Denardo, E.V., "Computing a bias-optimal policy in a discrete-time Markov decision problem," Opns. Res. 18, 279-289, 1970.

4 Ohno, K., "A unified approach to algorithms with a suboptimality λ^{test} in discounted semi-Markov decision processes," JORSJ 24, 296-324, 1981.