

On adaptive policies in Markovian Decision Processes  
with uncertain transition matrices

千葉大 教育 蔵野 正美

## 1. はじめに

推移確率行列が未知の場合のマルコフ決定過程の研究は色々な見地(たとえばゲーム論的, 最尤法, ベイズ的方法)から行われている。ここでは learning policy(学習政策)を定義して, discount case および average case において, learning policy と optimal policyとの関係を調べる。§.3では discount case におけるベイズ的方法が議論され, learning policy が定義される。そして, 任意の  $\epsilon > 0$  に対して,  $\epsilon$ -optimalな learning policy が存在することを証明する。また optimal policy は必ずしも learning policy でない例が与えられる。§.4では, average caseにおいて, optimalな learning policy が存在することを示す。

証明はいずれも構成的である。Sweat [10] は, a game model involving repeated play of a matrix game with unknown entries

において，“learning”と“optimal”との関係を研究した。この小論の手法は、いわば Sweat の方法の Markov model への適応と拡張といってよい。

## 2. Definitions and Notations

$S = \{1, 2, \dots, L\}$ : state space     $A = \{1, 2, \dots, K\}$ : action space

$r$ :  $S \times A$  の上に定義された return function

仮定     $0 \leq r(s, a) \leq R < \infty$

parameter space of  $K$  unknown stochastic matrices を

$$\Theta = \{ \theta / \theta = (\theta^a ; a \in A), \theta_{ss}^a \geq 0, \sum_{s \in S} \theta_{ss}^a = 1 \text{ for } s, s' \in S, a \in A \}$$

で表わす。可測空間  $(\Theta, \mathcal{F}_\Theta)$  の上の確率分布の全体を三で表わす。history  $w = (x_0, a_0, x_1, a_1, x_2, \dots)$  の全体を  $\Omega \equiv S \times (A \times S)^\infty$  で表わし、 $\mathcal{F}$  を  $\Omega$  の上の通常の  $\sigma$ -field とする。 $w \in \Omega$  に対して  $X_t(w) = x_t$  は  $t$  期の state を、 $A_t(w) = a_t$  は  $t$  期の action を、 $H_t(w) = (x_0, a_0, \dots, x_t)$  は  $t$  期までの history を表わす。但し、 $w = (x_0, a_0, \dots)$ 。

次に Policy を定義しよう。 $\pi = (\pi_0, \pi_1, \dots)$  が Policy とは  $t = 0, 1, 2, \dots$  において、 $\pi_t(\cdot / \xi, H_t(w))$  は  $\forall w \in \Omega, \forall \xi \in \Xi$  に対して  $A$  の上の確率分布であるときをいう。Policy の全体を  $\Pi$  で表わす。 $\forall \pi \in \Pi$  に対して、 $\pi_\xi = (\pi_{0,\xi}, \pi_{1,\xi}, \dots)$  を  $t = 0, 1, \dots$  に対して  $\pi_{t,\xi}(\cdot / H_t(w)) = \pi_t(\cdot / \xi, H_t(w))$  と定義する。 $\pi_\xi$  の全体を  $\Pi_\xi \equiv \{\pi_\xi ; \pi \in \Pi\}$  とする。

$\pi \in \Pi$ ,  $\theta \in \Theta$ ,  $\xi \in \Xi$  が与えられたときの system の推移法則は次の条件に従うものとする。

$$P(\Delta_t = a / X_0, \Delta_0, \dots, X_t = s) = \pi_{\xi, t}(a / X_0, \Delta_0, \dots, X_t = s)$$

$$P(X_{t+1} = s' / X_0, \dots, X_t = s, \Delta_t = a) = \theta_{ss'}^a$$

従って、通常の方法で初期値  $x \in S$  に対して  $(\Omega, \mathcal{F})$  の上の確率分布  $P_{\theta, \pi_{\xi}}^x(\cdot)$  が定義される。また  $\pi \in \Pi$ ,  $\xi \in \Xi$  が与えられたときの  $(\Omega, \mathcal{F})$  の上の確率分布  $P_{\xi, \pi}^x(\cdot)$  は次で定義する。  
 $P_{\xi, \pi}^x(D) = \int_{\Theta} \xi(d\theta) P_{\theta, \pi_{\xi}}^x(D) \quad \text{for } D \in \mathcal{F}$

### 3. Adaptive and learning policies in the discounted case.

$\beta$  を discount factor  $0 \leq \beta < 1$  として、 $\pi \in \Pi$ ,  $\xi \in \Xi$ ,  $x \in S$  に対して、無限期間の期待利得、最適期待利得を次のように定める。

$$V_{\beta}[\pi](x, \xi) \equiv E_{\xi, \pi}^x \left[ \sum_{t=0}^{\infty} \beta^t r(x_t, \Delta_t) \right]$$

$$V_{\beta}(x, \xi) \equiv \sup_{\pi \in \Pi} V_{\beta}[\pi](x, \xi)$$

#### 定義 1

$\pi^*$ ; optimal  $\Leftrightarrow V_{\beta}(x, \xi) = V_{\beta}[\pi^*](x, \xi)$  for  $x \in S$ ,  $\xi \in \Xi$ .

$\forall \varepsilon > 0$ ,  $\pi^*$ ;  $\varepsilon$ -optimal  $\Leftrightarrow V_{\beta}[\pi^*](x, \xi) \geq V_{\beta}(x, \xi) - \varepsilon$

for  $x \in S$ ,  $\xi \in \Xi$ .

$a \in A$ ,  $s, s' \in S$  に対して  $\Xi$  の上の operator  $T_{ss'}^a$  を次で定義する。

$$T_{ss'}^a \xi(D) = \int_D \theta_{ss'}^a d\xi(\theta) / \int_{\Theta} \theta_{ss'}^a d\xi(\theta) \quad \text{for } D \in \mathcal{F}_{\Theta} \quad (\text{Bayes の定理})$$

次の定理は chap 3 of Martin [7] につかわれたのと同じ方法で証明される。

### 定理1 ([7])

$\{v_\beta(x, \xi)\}$  は次の関数方程式を満足する。

$$(1) \quad v_\beta(x, \xi) = \max_{a \in A} \left\{ r(x, a) + \beta \sum_{x' \in S} \bar{\theta}_{xx'}^a(\xi) v_\beta(x', T_{xx'}^a(\xi)) \right\}$$

for  $x \in S$  and  $\xi \in \Xi$  但し,  $\bar{\theta}_{xx'}^a(\xi) = \int_{\Theta} \theta_{xx'}^a d\xi$ .

各  $x \in S$ ,  $\xi \in \Xi$  に対して  $g^*(x, \xi)$  を (1) の右辺を最大にする action の 1 つを表わすとする。そのとき,  $\pi^* = (\pi_0^*, \pi_1^*, \dots)$  を  $t=0, 1, 2, \dots$  に対して  $\pi_t^*(\{g^*(x_t, \xi_{H_t})\} / \xi, H_t) = 1$  と定義すると,  $\pi^*$  は optimal policy である。但し  $\xi_{H_t}$  は, history  $H_t$  が与えられたときの事後分布である。

$\pi \in \Pi$  が与えられたとき,  $\xi \in \Xi$  と  $n$  期までの history  $H_n(w)$  に対して policy  $\pi(\xi, H_n) = (\pi(\xi, H_n)_0, \pi(\xi, H_n)_1, \dots)$  を  $\pi(\xi, H_n)_t(\cdot / \xi', H_t(w')) = \pi(\cdot / \xi, H_n, H_t(w'))$  ( $t=0, 1, 2, \dots$ ) for  $w' \in \Omega$ ,  $\xi' \in \Xi$  と定義する。

$\Theta^+ \equiv \{\theta / \theta_{ss'}^a > 0 \text{ for } s, s' \in S \text{ and } a \in A\}$  として, learning policy を次に定義しよう。

### 定義2

$\delta > 0$  に対して,  $\pi$  が  $\delta$ -learning policy である  
 $\Leftrightarrow \forall \xi \in \Xi, \forall \theta \in \Theta^+ \cap \text{supp } \xi$  に対して次が成立する:  $\forall \varepsilon > 0$  に対して次の(i), (ii) を満足する  $A \in \mathbb{F}$  が存在する

- (i)  $P_{\theta, \pi_\xi}^x(A_\xi) = 1$  for  $x \in S$  and  
(ii)  $\forall w \in A_\xi$  に対して正の整数  $n_\xi(w)$  が存在して,  $\forall n \geq n_\xi(w)$   
 $\forall$  に対して  $\theta$  が parameter の実際の値のとき, policy  $\pi(\xi, H_n(w))$   
は  $(\delta + \varepsilon)$ -optimal となる。

i.e.  $V_\beta[\pi(\xi, H_n(w))] (x, I_\theta) \geq V_\beta(x, I_\theta) - (\varepsilon + \delta)$  for  $\forall x \in S$ ,  
 $\forall n \geq n_\xi(w)$ . ここで  $I_\theta$  は  $\theta$  で退化してある確率分布を表  
わす。

$\delta$ -learning policy を使うと parameter の実際の値に対して,  
system の operation は確率 1 で最終的に  $\delta$ -最適となる。 $\pi$  が  
任意の  $\delta > 0$  に対して  $\delta$ -learning ならば  $\pi$  を learning policy  
という。

## 定理 2

$\forall \varepsilon > 0$  に対して  $\varepsilon / (1-\beta)$ -optimal な learning policy が存  
在する。

以下, 定理 2 の証明のために準備をする。

$s \in S, a \in A$  に対して

$$\begin{aligned} u_t(s, a)(w) &= 1 \quad \text{if } X_t(w) = s \text{ and } \Delta_t(w) = a \\ &= 0 \quad \text{otherwise} \end{aligned}$$

$$W \equiv \{ w / \sum_{t=0}^{\infty} u_t(s, a) = +\infty \text{ for } s \in S, a \in A \}$$

そのとき次の lemma を得る。

Lemma 1

$\pi = (\pi_0, \pi_1, \dots)$  を次を満たす policy とする。すなわち、次の(i),(ii)を満足する正の実数列  $\{\varepsilon_k\}_{k=0}^{\infty}$  が存在する：

- (i)  $\pi_t(a/\xi, H_k(w)) \geq \varepsilon_k \quad \text{for } a \in A, \xi \in \Xi, w \in \Omega$
  - (ii)  $\sum_0^{\infty} \varepsilon_k = +\infty$
- $$\Rightarrow P_{\theta, \pi_\xi}^x(w) = 1 \quad \text{for } \theta \in \Theta^+, x \in S, \xi \in \Xi$$

<略証>

任意の  $(s, a) \in S \times A$  を固定する。簡単のために  $P_{\theta, \pi_\xi}^x(\cdot)$  を  $P(\cdot)$  で表わす。  
 $B_t \equiv \{u_t(s, a) = 1\} \quad (t=1, 2, \dots)$

任意の  $n, N$  ( $N > n$ ) に対して

$$P\left(\bigcap_{t=n}^N B_t^c\right) \leq P(B_n^c) e^{-\sum_{t=n+1}^N P(B_t^c / B_m^c \cap \dots \cap B_{t-1}^c)}$$

仮定(i)より  $P(B_t^c / B_m^c \cap \dots \cap B_{t-1}^c) \geq [\min_{s, s', a} \theta_{ss'}^a] \varepsilon_t$

$$\text{ゆえに } P\left(\bigcap_{t=n}^N B_t^c\right) \leq P(B_n^c) e^{-[\min_{s, s', a} \theta_{ss'}^a] \sum_{t=n+1}^N \varepsilon_t}$$

仮定(ii)より  $P(\lim_{t \rightarrow \infty} B_t^c) = 0$

すなわち  $P(\overline{\lim_{t \rightarrow \infty}} B_t) = 1$

QED.

Lemma 2 ([8], [9])

$$\xi \in \Xi, \theta \in \Theta^+ \cap \text{supp } \xi, x \in S, \pi \in \Pi$$

をとく

$$P_{\theta, \pi_\xi}^x(w) = 1 \Rightarrow \xi_{H_n(w)} \rightarrow I_\theta \text{ weakly in } \Xi$$

as  $n \rightarrow \infty$  with  $P_{\theta, \pi_\xi}^x$ -probability 1

Lemma 3

$\xi_n \rightarrow \xi$  weakly in  $\Xi$  as  $n \rightarrow \infty$   
 $\Rightarrow V_\beta(x, \xi_n) \rightarrow V_\beta(x, \xi)$  for  $x \in S$

<略証>

関数列  $\{U^l(x, \xi), l=0, 1, 2, \dots\}$  を逐次的に定義する。

$$U^l(x, \xi) = \max_{a \in A} \left\{ r(x, a) + \beta \sum_{x' \in S} \bar{\theta}_{xx'}^a(\xi) U^{l-1}(x', T_{xx'}^a(\xi)) \right\}$$

$$U^0(x, \xi) = 0$$

そのとき Helly-Bray theorem と induction によつて、

$\xi_n \rightarrow \xi$  weakly in  $\Xi$  as  $n \rightarrow \infty$  のとき、

$U^l(x, \xi_n) \rightarrow U^l(x, \xi)$  ( $l=1, 2, \dots$ ) が成り立つ。

従つて、 $\forall \varepsilon > 0$  と各  $l$  ( $l=1, 2, \dots$ ) に対して正の整数  $N_l$  が存在して  $|U^l(x, \xi_n) - U^l(x, \xi)| < \frac{\varepsilon}{3}$  for  $n \geq N_l$  が成り立つ。

また  $U^l(x, \xi) \rightarrow V_\beta(x, \xi)$  uniformly in  $(x, \xi)$  as  $l \rightarrow \infty$  が成り立つことはしめされるので、従つて、ある正の整数  $M$  が存在して、 $l \geq M$  に対して

$$|V_\beta(x, \xi_n) - U^l(x, \xi_n)| < \frac{\varepsilon}{3},$$

$$|V_\beta(x, \xi) - U^l(x, \xi)| < \frac{\varepsilon}{3} \quad \text{が成り立つ。}$$

結局、 $l \geq M$  に対して

$$|V_\beta(x, \xi_n) - V_\beta(x, \xi)| < \frac{\varepsilon}{3} + \frac{\varepsilon}{3} + \frac{\varepsilon}{3} = \varepsilon \quad \text{for } n \geq N_l$$

を得る。  $\square \text{ E.D.}$

## &lt;定理の略証&gt;

$\forall \varepsilon > 0$ ,  $\delta_\varepsilon \equiv \varepsilon(1-\beta) / R(L-1)$  とする。

policy  ${}^\varepsilon\pi = ({}^\varepsilon\pi_0, {}^\varepsilon\pi_1, \dots)$  を次のように定義する。

$\xi \in \Xi$ ,  $w \in \Omega$ ,  $t=0, 1, 2, \dots$  に対して

$${}^\varepsilon\pi_t(a/\xi, H_t(w)) = \begin{cases} 1 - (K-1)\delta_\varepsilon & \text{if } a = g^*(x_t(w), \xi_{H_t(w)}) \\ \delta_\varepsilon & \text{otherwise} \end{cases}$$

但し,  $g^*$  は定理 1 で定義された。

与えられた  $\varepsilon > 0$  に対して,  $\{\varepsilon_t\}_0^\infty$  を次をみた可数列とする。

$$(a) \quad \varepsilon_t \leq \varepsilon \quad \text{for } t=0, 1, 2, \dots$$

$$(b) \quad \varepsilon_t \downarrow 0 \quad \text{as } t \rightarrow \infty \quad \text{and} \quad (c) \quad \sum_0^\infty \varepsilon_t = +\infty$$

そのとき,  ${}^\varepsilon\pi^* = ({}^\varepsilon\pi_0^*, {}^\varepsilon\pi_1^*, \dots)$  を  $t=0, 1, 2, \dots$  に対して

$${}^\varepsilon\pi_t^*(a/\xi, H_t(w)) = {}^{\varepsilon_t}\pi_t(a/\xi, H_t(w)) \quad \text{for } w \in \Omega, \xi \in \Xi$$

で定義する。今,  ${}^\varepsilon\pi^*$  は  $\varepsilon/(1-\beta)$ -optimal かつ learning policy となることを示そう。 $M(S \times \Xi)$  を  $S \times \Xi$  の上に定義された有界な Bain 関数の全体とする。 $g \in M(S \times \Xi)$  によ,  $n$  期で truncate される過程を考えて

$$V_\beta^n[\pi, g(\cdot, \cdot)](x, \xi) = E_{\xi, \pi}^x \left[ \sum_0^n \beta^t r(x_t, \Delta_t) + \beta^{n+1} g(x_{n+1}, \xi_{H_{n+1}}) \right].$$

そのとき, 任意の  $n$  ( $n=1, 2, \dots$ ) に対して

$$V_\beta^n[{}^\varepsilon\pi^*, V_\beta(\cdot, \cdot) - \varepsilon/(1-\beta)](x, \xi) \leq V_\beta(x, \xi) - \frac{\varepsilon}{1-\beta} \quad \text{for } x \in S, \xi \in \Xi$$

が示される。また

$$n \rightarrow \infty \text{ のとき } V_\beta^n[{}^\varepsilon\pi^*, V_\beta(\cdot, \cdot) - \varepsilon/(1-\beta)](x, \xi) \rightarrow V_\beta[{}^\varepsilon\pi^*](x, \xi)$$

であるから、 $\varepsilon\pi^*$  は  $\varepsilon/(1-\beta)$ -optimal となる。

次に証明は複雑になるが、lemma 1 ~ 3 および truncation method によって  $\varepsilon\pi^*$  が learning policy であることが示される。

$$Q \in D.$$

例 (optimal な learning policy は必ずしも存在しない)

$$S = A = \{1, 2\}, r(1, a) = 1 \text{ for all } a \in A,$$

$$r(2, a) = 0 \text{ for all } a \in A$$

$$\bar{\Theta} \equiv \left\{ \theta / \theta'_{21} = \theta'_{22} = \frac{1}{2}, \theta_{11}^a = \theta_{12}^a = \frac{1}{2} \text{ for all } a \in A \right\}$$

$$\Xi^* = \left\{ \xi \in \Xi / \text{supp } \xi = \bar{\Theta} \right\}$$

定理 1 から  $\forall \xi \in \Xi^*$  に対して、次式を得る。

$$\begin{aligned} V_\beta(2, \xi) &= \max \left\{ \frac{\beta}{2} (V_\beta(1, T'_{21}\xi) + V_\beta(2, T'_{22}\xi)), \right. \\ &\quad \left. \beta (\bar{\theta}_{21}^2(\xi) V_\beta(1, T_{21}^2\xi) + (1 - \bar{\theta}_{21}^2(\xi)) V_\beta(2, T_{21}^2\xi)) \right\} \end{aligned}$$

$\theta_{21}^{*2} = \frac{1}{9}$  なる  $\theta^* \in \bar{\Theta}$  を選ぶ。このとき

$\exists \{\xi_n\}_1^\infty$  such that

(a)  $\xi_n \rightarrow I_{\theta^*}$  weakly in  $\Xi$  as  $n \rightarrow \infty$

(b)  $\xi_n \in \Xi^*$  for  $n = 1, 2, 3, \dots$

lemma 3 & Helly-Bray theorem から

$$\begin{aligned} V_\beta(2, \xi_n) &\rightarrow V_\beta(2, I_{\theta^*}) = \max \left\{ \frac{1}{2} \beta (V_\beta(1, I_{\theta^*}) + V_\beta(2, I_{\theta^*})) \right. \\ &\quad \left. , \beta \left( \frac{1}{9} V_\beta(1, I_{\theta^*}) + \frac{8}{9} V_\beta(2, I_{\theta^*}) \right) \right\} \end{aligned}$$

$V_\beta(1, I_{\theta^*}) > V_\beta(2, I_{\theta^*})$  から

$$V_\beta(2, \xi^*) = \frac{1}{2} \beta (V_\beta(1, T'_{21}\xi^*) + V_\beta(2, T'_{22}\xi^*))$$

なる  $\xi^*$  が存在する。その結果  $P_{\xi^*, \pi^*}^2(\Delta_t=1) = 1$  for  $t=0, 1, \dots$  となる。すなわち, a prior 分布が  $\xi^*$  のとき optimal policy は常に action 1 のみを選択することになる。ところが, たとえば  $\theta_{21}^2 = \frac{8}{9}$  なる  $\theta \in \bar{\Theta}$  が実際の値のときは,  $\pi^*(\xi^*, H_n)$  は決して optimal にならない。

### Corollary

$P_{\theta, \pi_\xi^*}(W) = 1$  for  $\xi \in \Xi$ ,  $\theta \in \Theta^+$ ,  $x \in S$  なる optimal policy  $\pi^*$  が存在するならば,  $\pi^*$  はまた learning policy である。

### 4. Non-Bayesian approaches to MDP's with uncertain transition matrices.

この節では未知の推移確率 (parameter) に対する事前分布が存在しなく, かつ無限期間での一期あたりの平均期待利得を最大にする問題をとり扱い, 最尤法をもちいて optimal な learning policy を構成する。事前分布が存在する場合は Rose, J.S. [9] によって研究された。 $t=0, 1, 2, \dots$  に対して  $\gamma_t(\cdot | H_t(w))$  は  $A$  の上の確率分布であるとき  $\gamma = (\gamma_0, \gamma_1, \gamma_2, \dots)$  をこの節では policy と呼ぶ。policy の全体を  $\Gamma$  で表わす。さらに,  $\theta$  が真のとき初期値  $x \in S$  で policy  $\gamma$  をつかったときの  $(\Omega, \mathcal{F})$  の上の確率分布を  $P_{\theta, \gamma}^x(\cdot)$  で表わす。

$x \in S$ ,  $\gamma \in \Gamma$ ,  $\theta \in \Theta$  に対して T 期までの総期待利得を

$\psi_T[\gamma](x, \theta) \equiv E_{\theta, \gamma}^x \left[ \sum_0^{T-1} r(x_t, \Delta_t) \right]$  として、無限期間における平均期待利得と最適平均期待利得を次のように定義する：

$$\psi[\gamma](x, \theta) \equiv \lim_{T \rightarrow \infty} \frac{1}{T} \psi_T[\gamma](x, \theta)$$

$$\psi(x, \theta) = \sup_{\gamma \in P} \psi[\gamma](x, \theta)$$

### 定義 3

$\gamma^*$ ; optimal in the average case

$$\Leftrightarrow \psi(x, \theta) = \psi[\gamma^*](x, \theta) \quad \text{for } x \in S, \theta \in \Theta^+$$

$\gamma = (\gamma_0, \gamma_1, \dots)$ ; Markov

$$\Leftrightarrow \gamma_t(\cdot / x_0, \dots, x_t) = \gamma_t(\cdot / x_t) \quad t=0, 1, 2, \dots$$

$\gamma = (\gamma_0, \gamma_1, \dots)$ ; stationary

$$\Leftrightarrow \gamma_t(\cdot / x_0, \Delta_0, \dots, x_t) = \gamma_0(\cdot / x_t) \quad t=0, 1, 2, \dots$$

$\gamma$  が stationary のとき、 $\gamma$  を  $\gamma_0^{(\infty)}$  と記す。また、deterministic stationary policy を  $f^\infty = (f, f, \dots)$  と表わす。

### Lemma 1 ([2])

$\theta \in \Theta^+$  が既知とする。そのとき、optimal deterministic policy  $f_\theta^{(\infty)}$  が存在する。また、 $\psi[f_\theta^{(\infty)}](x, \theta)$  は初期値  $x$  に無関係である。

そこで、 $\theta \in \Theta^+$  に対して、 $g_\theta \equiv \psi[f_\theta^{(\infty)}](x, \theta)$  とする。

任意の stationary policy  $\gamma^{(\infty)} = (\gamma, \gamma, \dots)$  と  $\theta \in \Theta$  に対して  $L \times L$  matrix  $\theta^\gamma = (\theta_{ss'}^\gamma)$  と limiting matrix  ${}^*\theta^\gamma$  を次のよう 定義する。

$$\theta_{ss'}^r = \sum_{a \in A} \theta_{ss'}^{a,r} r(a/s) \quad * \theta^r \equiv \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{t=0}^n \theta^{r(t)}$$

そのとき次のlemmaが成り立つことが知られている。

Lemma 2 ([I], [II])

$$\theta \in \Theta^+, \gamma^{(\infty)} = (r, r, \dots), \quad H^r \equiv [I - \theta^r + * \theta^r]^{-1} - * \theta^r$$

$$y^r = (y_1^r, \dots, y_L^r)' = H^r r(r)$$

$\Rightarrow$

$$(i) \quad \psi[\gamma^{(\infty)}](s, \theta) = \sum_{s' \in S} * \theta_{ss'}^r r_{s'}(r) \quad \text{and} \quad \sum_{s' \in S} * \theta_{ss'}^r y_{s'}^r = 0$$

$$(ii) \quad r_s(r) + \sum_{s' \in S} \theta_{ss'}^r y_{s'}^r = \psi[\gamma^{(\infty)}](s, \theta) + y_s^r$$

$$(iii) \quad \psi_T[\gamma^{(\infty)}](s, \theta) = T \psi[\gamma^{(\infty)}](s, \theta) + y_s^r - \sum_{s' \in S} \theta_{ss'}^{r(T)} y_{s'}^r$$

但し、L-vector  $r(r) = (r_1(r), \dots, r_L(r))'$  は

$$r_s(r) = \sum_{a \in A} r(s, a) r(a/s) \quad s \in S \quad \text{で定義される。}$$

次に  $^k \varepsilon > 0$  ( $0 < \varepsilon < \frac{1}{L}$ ),  $^k \theta \in \Theta^+$  に対して, stationary policy  ${}^k \gamma_{\theta}^{(\infty)}$  を次のようく定義する。

$${}^k \gamma_{\theta}(a/s) = \begin{cases} 1 - (L-1)\varepsilon & \text{if } a = f_{\theta}(s) \\ \varepsilon & \text{otherwise} \end{cases}$$

そのとき、次のlemmaを得る。

Lemma 3

$$\theta(n) \in \Theta^+ \rightarrow \theta \in \Theta^+ \text{ as } n \rightarrow \infty$$

$$\Rightarrow \psi[{}^k \gamma_{\theta(n)}^{(\infty)}](x, \theta) \rightarrow g_{\theta} \quad \text{as } n \rightarrow \infty \text{ and } \varepsilon \rightarrow 0$$

<略証>

記号を簡単にするために  $f_n = f_{\theta(n)}$ ,  ${}^k \gamma_n = {}^k \gamma_{\theta(n)}$  とおく。

Denard & Miller [3] より、各  $n$  ( $n=1, 2, \dots$ ) に対して、次  
の(1), (2)を満足する L-vector  $\tilde{y}^n = (\tilde{y}_1^n, \dots, \tilde{y}_L^n)$  が存在する。

$$(1) \quad \Gamma(s, f_n(s)) + \sum_{s' \in S} \theta_{ss'}^{f_n(s)}(n) \tilde{y}_{s'}^n = g_{\theta(n)} + \tilde{y}_s^n$$

$$(2) \quad \Gamma(s, a) + \sum_{s' \in S} \theta_{ss'}^a(n) \tilde{y}_{s'}^n \leq g_{\theta(n)} + \tilde{y}_s^n \quad \text{for } a \in A$$

と定義されば、(2)から  
次を得る。

$$(3) \quad \Gamma(s, f_\theta(s)) + \sum_{s' \in S} \theta_{ss'}^{f_\theta(s)} \tilde{y}_{s'}^n + \sum_{s' \in S} Y_{ss'}^n(f_\theta) \tilde{y}_{s'}^n \leq g_{\theta(n)} + \tilde{y}_s^n$$

$\{\theta_{xs}^{f_\theta(x)}\}$  を  $\{\theta_{xs}^x\}$  の limiting matrix とすると、lemma 2 の(i)  
と(3)より次を得る。

$$(4) \quad g_\theta \leq g_{\theta(n)} - \sum_{s, s'} * \theta_{xs}^{f_\theta(x)} Y_{ss'}^n(f_\theta) \tilde{y}_{s'}^n$$

一方、 ${}^\varepsilon \gamma_n$  の定義から

$$\theta_{ss'}^{f_n(s)} = \theta_{ss'}^{\varepsilon \gamma_n} - \varepsilon \sum_{a \in A} (\theta_{ss'}^a - \theta_{ss'}^{f_n(s)})$$

$$\Gamma(s, f_n(s)) = \Gamma_s({}^\varepsilon \gamma_n) - \varepsilon \sum_{a \in A} (\Gamma(s, a) - \Gamma(s, f_n(s)))$$

$$\text{但し, } \Gamma_s({}^\varepsilon \gamma_n) = \sum_{a \in A} \Gamma(s, a) {}^\varepsilon \gamma_n(a)$$

これを用いて計算すると(1)から次を得る。

$$(5) \quad \psi[{}^\varepsilon \gamma_n^{(\infty)}](x, \theta) = g_{\theta(n)} + \varepsilon \sum_{s \in S} * \theta_{xs}^{\varepsilon \gamma_n} Z_s^n$$

$$\text{但し, } Z_s^n = \sum_{s', a} (\theta_{ss'}^a - \theta_{ss'}^{f_n(s)}) \tilde{y}_{s'}^n + \sum_a (\Gamma(s, a) - \Gamma(s, f_n(s))) \\ + \sum_{s'} Y_{ss'}^n(f_n) \tilde{y}_{s'}^n$$

(4) と(5)より次を得る。

$$g_\theta - \psi[{}^\varepsilon \gamma_n^{(\infty)}](x, \theta) \leq l_n \|\tilde{y}\| + \varepsilon (2R + 2\|\tilde{y}\| + l_n \|\tilde{y}\|)$$

$$\text{但し, } l_n = \max_{x, x' \in S, a \in A} |\theta_{xx'}^a(n) - \theta_{xx'}^a|, \quad \|\tilde{y}\| = \left( \sum_{e=1}^L (\tilde{y}_e^n)^2 \right)^{\frac{1}{2}}$$

$\{y^n\}$  は有界な vector であるから,  $n \rightarrow \infty, \varepsilon \rightarrow 0$  のとき

$$g_0 - \psi[\varepsilon r_n^{(\infty)}](x, \theta) \rightarrow 0 \quad \text{となる。} \quad \text{QED.}$$

$s, s' \in S, a \in A$  に対して  $u_t(s, a, s')$  を次のように定義する。

$$u_t(s, a, s') = 1 \quad \text{if } X_t(w) = s, \Delta_t(w) = a, X_{t+1}(w) = s'$$

$$u_t(s, a, s') = 0 \quad \text{otherwise}$$

$\{N_\ell\}_0^\infty$  を  $N_0 = 0$  なる正の整数の数列とする。

$s, s' \in S, a \in A, k=1, 2, \dots$  に対して,  $\hat{\theta}_{ss'}^a(k)$  を次のように定義する。

$$\hat{\theta}_{ss'}^a(k) = \frac{\sum_{\ell=1}^{L_k} u_\ell(s, a, s')}{\sum_{\ell=1}^{L_k} u_\ell(s, a)} \quad \text{但し, } L_k = \sum_{\ell=0}^k N_\ell.$$

matrix  $\hat{\theta}(k) = \{\hat{\theta}_{ss'}^a(k)\}$  は  $L_k$  期の  $\theta$  の最尤推定量である。

$\{\lambda_\ell\}_0^\infty$  を  $0 < \lambda_\ell < 1, \lambda_0 = 1$  なる正の実数列として,  ${}^A\hat{\theta}_0 \in \mathbb{H}^+$  に対して

$${}^*\hat{\theta}_k \equiv \lambda_k \hat{\theta}_0 + (1 - \lambda_k) \hat{\theta}(k), \quad k = 1, 2, \dots \quad \text{とする。}$$

次に,  $\{\varepsilon_\ell\}_0^\infty$  を  $0 < \varepsilon_\ell < \frac{1}{L}$  なる正の実数列として

policy  $\gamma^* = (\gamma_0^*, \gamma_1^*, \dots)$  を次のように定義する:

$$\gamma_t^*(\cdot / H_t(w)) = {}^{\varepsilon_k} \gamma_{*\hat{\theta}_k}(\cdot / X_t(w)) \quad \text{if } L_k \leq t < L_k + N_k = (L_{k+1})$$

そのとき, policy  $\gamma^*$  はある条件のもとで optimal & learning policy であることが次の定理で示される。

### 定理 3

$$(i) \quad \sum_1^\infty \varepsilon_\ell N_\ell = +\infty$$

(ii)  $N_\theta \uparrow +\infty$ ,  $\lambda_\theta \downarrow 0$ ,  $\varepsilon_\theta \downarrow 0$

$\Rightarrow \hat{\theta}_\theta \in \Theta^+$  に対して  $r^*$  は optimal である。

さらに,  $\psi[r^*](x, \theta) = g_\theta$  for all  $\theta \in \Theta^+$

<略証>

§3 の lemma 1 と大数の強法則から,  $k \rightarrow \infty$  のとき

${}^*\hat{\theta}_k \rightarrow \theta$  with  $P_{\theta, r^*}^x$ -probability 1 である。今, 簡単のため

を  $r_k \equiv {}^{\varepsilon_k}r^*_{\hat{\theta}_k}$  とする。そのとき

$\eta({}^*\hat{\theta}_k) \equiv \psi[r_k^{(0)}](x, \theta) - g_\theta$  すると, lemma 3 より,

$k \rightarrow \infty$  のとき  $\eta({}^*\hat{\theta}_k) \rightarrow 0$  with  $P_{\theta, r^*}$ -probability 1.

lemma 2 から  $k=0, 1, \dots$  に対して L-vector  $y^k = (y_1^k, \dots, y_L^k)$

が存在して, 次式が成り立つ。

$$\begin{aligned} & \sum_{t=L_k}^{L_{k+1}-1} E_{\theta, r_k^{(0)}}^x [r(x_t, \Delta_t) / X_{L_k} = x] \\ &= [g_\theta - \eta({}^*\hat{\theta}_k)] N_{k+1} + y_x^k - \sum_{x' \in S} \theta_{xx'}^{r_k} y_{x'}^k \end{aligned}$$

ゆえに次を得る。

$$\begin{aligned} \frac{1}{L_{k+1}} \sum_{t=0}^{L_{k+1}-1} E_{\theta, r^*}^x [r(x_t, \Delta_t)] &\geq g_\theta - \frac{1}{L_{k+1}} \sum_{k=0}^K E_{\theta, r^*}^x [\eta({}^*\hat{\theta}_k)] N_{k+1} \\ &\quad - 2 \frac{1}{L_{k+1}} \sum_{k=0}^K \|y^k\|. \end{aligned}$$

従って, bounded convergence theorem, Toeplitz Lemma [6]

等をつかって,

$$\begin{aligned} \psi[r^*](x, \theta) &= \lim_{K \rightarrow \infty} \frac{1}{L_{K+1}} \sum_{k=0}^{L_{K+1}-1} E_{\theta, r^*}^x [r(x_t, \Delta_t)] \\ &\geq g_\theta \quad \text{for all } \theta \in \Theta^+ \end{aligned}$$

一方,  $\psi[r](x, \theta) \leq g_\theta$  for  $r \in \Gamma$  より

$\psi[r^*](x, \theta) = g_\theta$  for all  $\theta \in \Theta^+$  が成り立つ。

QED.

### REFERENCE

[1] Blackwell, D. (1962)

Discrete Dynamic programming.

Ann. Math. Statist. 33, 719-726.

[2] Derman, Cyrus (1962)

On sequential decision and Markov chains.

Mag. Sci. 9, 16-24.

[3] Denardo, E. V. and Miller, B. L. (1968)

An optimality condition for discrete dynamic programming  
with no discounting.

Ann. Math. Statist. 39, 1220-1227.

[4] Howard, Ronald A. (1960)

Dynamic programming and Markov processes.

the M. I. T. press.

[5] M. Kurano (1972)

Discrete-time Markovian decision processes with an unknown  
parameter, -average return criterion-

- Journal of Op. Reser. of Japan. 15 no. 2 67-76.
- [6] Loeve, M. (1960)  
Probability Theory, Second Edition.  
 D. Van Nostrand Co., Inc., New Jersey.
- [7] Martin, J. J. (1967)  
Bayesian Decision problem and Markov chains.  
 Wiley, New York.
- [8] Rose, J. S. (1971)  
 Markovian Decision Processes under uncertainty.  
 PH. D. dissertation, Northwestern U.
- [9] Rose, J. S. (1975)  
 Markov Decision Processes under uncertainty - average  
 return criterion. Unpublished report.
- [10] C. W. Sweat (1968)  
 Adaptive competitive decision in repeated play of a matrix  
 game with uncertain entries.  
 Naval R. L. Q. 15, 425-448.
- [11] Veinott, A. F. (1966)  
 On finding optimal policies in discrete dynamic programming  
 with no discounting.  
 Ann. Math. Stat., 37, 1284-1294.