

IV 決定過程

最適停止とともに決定過程

九大 理 古川長太

§1. 序

離散時間確率的系に対する最適化問題では、stage数を
予め指定せずに無限 stage にわたっての政策を問題とする
ならば、通常、評価関数として無限和 (stageについての
和)、またはその 1 stage 当りの平均の型を採用する。こ
こでは無限和の型の評価関数を最大化することに限って
考えることにする。

無限和の評価関数を採用することは、一応、我々が無限
stageにわたって行動(続けることを前提としているが)現
実にはいつかは stop せざるを得ない。また、もしろ stop し
た方が有利な場合もある。また、いわゆる sequential
analysis のように terminal decision が要求
されるときは、当然 いつかは stop しなくては 問題の

解決にならぬ。

ここに一つの決定過程 (Π) が与えられたとし、 Π の状態空間を S 、行動空間を A とする。 S に仮想的な状態 s_0 (吸收点) を加え、 $A_0 = \text{stop}$ を表わす行動 a_0 を加えて、 $S' = S + \{s_0\}$, $A' = A + \{a_0\}$ を新たにそれぞれ状態空間、行動空間として新しい決定過程 (Π') を構成すると、 Π' は stop するという行動をも含んだ決定過程になる。 Π' に対して無限和の評価関数を採用すると、表面上は無限和であるが " a_0 をとった時点" では process は実際には stop し、したがって、はじめに述べた問題は、このような定式化により通常の最適化問題へ帰着されるかのように見える。しかしながら、 Π' における最適政策が Π において果して確率 1 で stop するなどになるかどうか分らないし、また Π で stop する確率を計算することも一般に仲々困難である。

以上の理由により、ここに全く新しい定式化を考之なおすべくする。この報告では "確率 1 で stop する政策" (非常に rough を表現であるが) を考之、そのような政策の集合の中で最適なものを追究する。正確な定義は次の章にゆずるとして、さらに詳しく述べれば、通常の意味の政策 π と、停止時間 τ を pair として考之、あるゆる pair (π, τ) の集合の中で与えられた評価関数を最大化

することを問題とする。停止時間は停止規則と同等であるから、pariの政策(π, \bar{c})はまた、通常の意味の政策と停止則則とのpariともみなせる。

近來、ソ連のマルコフ過程論(特に強マルコフ過程)の発展に伴って、E.B. Dynkin および A.N. Shiryaev 等がいくつかの具体的な最適停止問題について、興味ある結果を与えた。またアメリカでは古くは Y.S. Chow, H. Robbins, 近年では L.E. Dubins, L.J. Savage 等により多くの具体的な最適停止問題が扱われて来た。

この報告は、定式化としては、前述のように通常の stochastic control problem に停止規則を入れた型でもあるし、また、別の面から見ると、上記の最適停止問題の一般化にさらに通常の control action を入れた型ともみなせる。

この報告は全体を通じて、上記の意味での pari の政策の class における最適政策(最適性の標準は今までいくつか与えられる)の存在定理と、最適政策の性質について論じたものである。その他、これらに付隨した多くの定性的な興味ある結果が与えられている。最適解の構成に関しては、状態空間、行動空間をごく簡単なものに限定すれば容易であるが、この報告でとり上げたような一般的な距離空間

では、可成り困難な問題で、この報告でも explicit な型で述べられ所までの追究はされていない。（存在定理の証明方法が有限 stage による政策改良の線に沿つているから、この意味では可成り constructive である）

以上のように、最適政策の構成法については不十分であるが、存在定理および最適政策の性質等については、一歩初期の目的を達したので、この報告を総括する。

§2. 記号および定式化

2.1 一般的な定義

Borel set ; complete separable metric space の Borel subset
 $(X, Y \text{ 等で表わす})$

Borel set X 上の probability measure ; X の Borel field 上の
probability measure

$P(X)$; X 上のすべての probability measure からなる空間

$\gamma(y|x)$; x を与えたときの Y 上の conditional prob. measure γ^x

(i) 各 $x \in X$ に対し $\gamma(\cdot|x)$ は Y 上の prob. measure

(ii) 各 Borel subset $B \subset Y$ に対し $\gamma(B|\cdot)$ は X 上の Baire 测度

$Q(Y|X)$; 定義空間 X, Y をきみて 上のようなすべての
conditional prob. measure からなる空間

$M(X)$; (i) §2 ~ §4 では X 上のすべての 実数値 Baire
関数からなる空間

(ii) §5, §6 では X 上のすべての 有界 Baire 関数
からなる空間

$u, v \in M(X)$ に対して $u \geq v$; $u(v) \geq v(x)$ for all $x \in X$

$p \in P(X)$, $u \in M(X)$ に対して pu ; $pu = \int_X u(v) dp(x)$

$g \in Q(Y|X)$, $u \in M(XY)$ に対して gu ; $gu(v) = \int_Y u(x, y) dg(y|x)$

$p \in P(X)$ が "degenerate"; $p\{x_0\} = 1$ for some $x_0 \in X$

$g \in Q(Y|X)$ が "degenerate"; $g(\cdot|x)$ が 各 x に対して degenerate

$g \in Q(Y|X)$ が "degenerate" であれば $g(\{f(x)\}|x) = 1$ for
all $x \in X$ なる X から Y への Baire 関数 f が存在する
ことになり、従って $u \in M(XY)$ なる u に対しては

$$f(u(v)) = u(x, f(v)) \quad \text{for all } x \in X$$

となる。

2.2 決定過程 κ 固有な定義

停止規則をともなう決定過程に対する最適化問題は
つきの 6 箇の 要素 S, A, g, r, g, α で決定される。

S ; 状態空間 (state space) を表わし, ある空でない Borel set

A ; 行動空間 (action space) を表わし, ある空でない Borel set

$\gamma = \{\gamma_1, \gamma_2, \dots\}$; system の運動法則を表わし, 各 n につき

$$\gamma_n \in Q(S|H_n A), H_n = S A S A \cdots S \quad (2n-1)$$

$r \in M(SAS)$; 利得関数 (reward function)

g ; 最終利得関数で S 上の有界 Baire 関数

α ; 割引率 (discount factor)

$$(i) \quad \$2 \sim \$4 \text{ では } 0 \leq \alpha \leq 1$$

$$(ii) \quad \$5, \$6 \text{ では } 0 \leq \alpha < 1$$

→ ここで, policy, stopping time より stopped policy 等についての定義を述べる。

政策 (policy) $\pi = \{\pi_1, \pi_2, \dots\}$; $\pi_n \in Q(A|H_n)$ $n=1, 2, \dots$

Markov policy; $\pi = \{\pi_1, \pi_2, \dots\}$ において各 π_n が $Q(A|S)$
の degenerate の要素

Markov policy を $\{f_1, f_2, \dots\}$ (ただし $f_n : S$ から
 A への Baire 関数) と書く。

定常政策 (stationary policy); $\pi = \{f_1, f_2, \dots\}$ が

Markov policy かつ $f_n = f$ for all n

stationary policy を f^∞ と書く。

$\pi = \{\pi_1, \pi_2, \dots\}$ に対し ${}^n\pi = \{\pi_{n+1}, \pi_{n+2}, \dots\}$ と
書く。故に ${}^0\pi = \pi$

π を任意の policy とすれば、 π と g_1, g_2, \dots つきの ような conditional probability measure e_π が定義される。

$$e_\pi = \pi_1 g_1 \pi_2 g_2 \dots \in Q(\text{ASAS} \dots | S)$$

\mathcal{G} ; S の σ -field

\mathcal{C} ; A の σ -field

標本空間 (sample space); $\Omega = SASA \dots$

標本点 (sample point) 又は 道 (path); $\omega = (s_1, a_1, s_2, a_2, \dots)$

$$\mathcal{F}_n; \{ \omega \mid s_i \in E_i, a_i \in F_i, s_2 \in E_2, a_2 \in F_3, \dots, s_n \in E_n \}$$

$$(ただし \quad E_i \in \mathcal{G}, F_i \in \mathcal{C}, i=1, 2, \dots, n)$$

の型の ω -集合のあらゆる和の族

(\mathcal{F}_n は 明らかに σ -field になる)

policy π に付随する停止時間 (stopping time)

$$t = t(\omega); \begin{cases} \text{(i) 正整数値の確率度数} \\ \text{(ii) } e_\pi[\{t(\omega) < \infty\}_{s_1} \mid s_1] = 1 \text{ for all } s_1 \in S \\ \text{(iii) } \{t(\omega) = n\} \in \mathcal{F}_n \text{ for each } n \end{cases}$$

$C(\pi)$; π に付随するすべての stopping time の集合

停止政策 (stopped-policy) (π, t) ;

π は 任意の policy, $t \in C(\pi)$ のとき (π, t) を stopped-policy といい、簡単のために s -policy とかく

$C^N(\pi)$; $C(\pi)$ に属つかう

$$e_{\pi}[\{\tau(\omega) \leq N\}_{s_1} | s_1] = 1 \quad \text{for all } s_1 \in S$$

をみたすすべての stopping time の集合

切り政策 (truncated-policy) ; $s\text{-policy } (\pi, t)$ で
あってかつ $t \in C^N(\pi)$ なる正整数 N が存在するもの
を truncated-policy といい、簡単のために $t\text{-policy}$
とかく

B_n ; $\{\omega | s_1 \in E_1, s_2 \in E_2, \dots, s_n \in E_n\}$ ($\bigcap_{E_i \in \mathcal{G}} i=1, 2, \dots, n$)
の型の ω -集合のあらゆる和の族
(B_n は明らかに σ -field)

\mathcal{G}_n ; $\{\omega | s_n \in E\}$ ($\bigcap E \in \mathcal{G}$)
の型のあらゆる ω -集合の族

$t \in C(\pi)$ が "Markov stopping time" ;

- (i) $\{\tau(\omega) = n\} \in B_n$ for each n
- (ii) $\{\tau(\omega) = n\} = \{\tau(\omega) > n-1\} \cap \Delta_n$
たゞ $\Delta_n \in \mathcal{G}_n$, for each n

Markov $s\text{-policy } (\pi, t)$; $\begin{cases} \text{(i) } \pi \text{ が "Markov policy"} \\ \text{(ii) } t \in C(\pi) \text{ かつ } t \text{ が "Markov stopping time"} \end{cases}$

stationary $t\text{-policy } (\pi, t)$; (π, t) が "Markov $s\text{-policy}$ " かつ, π が "stationary policy"

$t \in C(\pi)$ が stationary stopping time ; 大が "Markov" かつ Δ_n が n に無関数

つぎに評価関数の設定と、いくつかの最適性の規準を与える。

$$x_n = \sum_{k=1}^{n-1} \alpha^{k-1} r_k + \alpha^{n-1} g_n, \quad n=1, 2, \dots$$

これを簡単のために

$$x_n = \sum_{k=1}^{n-1} \alpha^{k-1} r_k + \alpha^{n-1} g_n, \quad n=1, 2, \dots$$

とかくこともある。

E^π ; e_π に関する expectation を表す積分作用素
 π -policy (π, t) から得られた総期待利得を評価関数として採用する。すなわち

$$\begin{aligned} E^\pi(x_t) &= E^\pi \left[\sum_{k=1}^{t-1} \alpha^{k-1} r_k + \alpha^{t-1} g_t \right] \\ &= e_\pi \left[\sum_{k=1}^{t-1} \alpha^{k-1} r_k + \alpha^{t-1} g_t \right] \end{aligned}$$

我々の目的は、上に定義した $E^\pi(x_t)$ を、何らかの意味で (π, t) について最大化することである。

Π^N ; N番 stage までに対するすべての policy $\{\pi_1, \pi_2, \dots, \pi_N\}$ の集合

$$\Lambda^N \equiv \{(\pi, t) \mid \pi \in \Pi^{N-1}, t \in C^N(\pi)\}$$

Π ; 無限 stage にわたってのすべての policy の集合

$$\Lambda \equiv \{(\pi, t) \mid \pi \in \Pi, t \in C(\pi)\}$$

$(\rho, \varepsilon)^N$ -optimal ; $(\hat{\pi}, \hat{t}) \in \Lambda^N$ かつ

$$\Pr\{E^{\hat{\pi}}(x_{\hat{t}}) \geq E^{\pi}(x_t) - \varepsilon\} = 1 \text{ for } \forall (\pi, t) \in \Lambda^N$$

のとき $(\hat{\pi}, \hat{t})$ を $(\rho, \varepsilon)^N$ -optimal とする

$(\rho, \varepsilon, \delta)$ -optimal ; $(\hat{\pi}, \hat{t}) \in \Lambda$ かつ

$$\Pr\{E^{\hat{\pi}}(x_{\hat{t}}) \geq E^{\pi}(x_t) - \varepsilon\} \geq 1 - \delta \text{ for } \forall (\pi, t) \in \Lambda$$

のとき $(\hat{\pi}, \hat{t})$ を $(\rho, \varepsilon, \delta)$ -optimal とする

(ρ, ε) -optimal ; $(\hat{\pi}, \hat{t})$ が $(\rho, \varepsilon, 0)$ -optimal のとき, これ
を (ρ, ε) -optimal とする

ε -optimal ; $(\hat{\pi}, \hat{t}) \in \Lambda$ かつ

$$E^{\hat{\pi}}(x_{\hat{t}}) \geq E^{\pi}(x_t) - \varepsilon \text{ for } \forall (\pi, t) \in \Lambda$$

のとき $(\hat{\pi}, \hat{t})$ を ε -optimal とする

optimal ; $(\hat{\pi}, \hat{t})$ が 0-optimal のときこれを optimal
とする

(ρ, ε) -dominate ; $\Pr\{E^{\hat{\pi}}(x_{\hat{t}}) \geq E^{\tilde{\pi}}(x_{\tilde{t}}) - \varepsilon\} = 1$ のとき

$(\hat{\pi}, \hat{t})$ は $(\tilde{\pi}, \tilde{t})$ を (ρ, ε) -dominate すと

とする

§ 3. $(\rho, \varepsilon, \delta)$ -optimal policy および (ρ, ε) -optimal policy の
存在

先づ本論に入る前に $\{x_m\}$ -process から $\{\beta_m^N(\pi)\}$ -process

を構成し それにもとづく 2 つの proposition をあげる。

π を任意の policy として 各 $N \geq 1$ に対し, $\{\beta_n^N(\pi)\}$ -process を 次式で backward に 定義する。

$$\left\{ \begin{array}{l} \beta_n^N(\pi) = \max \left[x_n, E^{n-1\pi}(\beta_{n+1}^N(\pi)) \right], \quad n=1, 2, \dots, N-1 \\ \beta_N^N(\pi) = x_N \end{array} \right. \quad (3.1)$$

ただし 各 x_n は E^π に 関して 可積分であるとし, $E^{n-1\pi}$ は conditional prob. measure $\pi_n \otimes \pi_{n+1} \otimes \pi_{n+2} \dots$ に 関する 積分を表わす。つたがって policy に関する 仮定から, 各 N , n につき $\beta_n^N(\pi)$ は \mathcal{F}_n -可測である。

$\{\beta_n^N(\pi)\}$ を用いて τ_N を 定義する。

$$\tau_N(\pi) = \text{the first } n \text{ such that } \beta_n^N(\pi) = x_n. \quad (3.2)$$

(3.1) により $\beta_N^N(\pi) = x_N$ だから $\tau_N(\pi) \in C^N(\pi)$.

今后簡単のために、誤解の恐れのない限り $\tau_N(\pi)$ を τ_N と 表わす。 τ_N の 定義から 明らかに

$$x_{\tau_N} = \beta_{\tau_N}^N(\pi).$$

Proposition 1 (3.1), (3.2) で 定義された $\{\beta_n^N(\pi)\}$, τ_N に対して

$$(a) \quad E^\pi(x_{\tau_N}) = \sup_{t \in C^N(\pi)} E^\pi(x_t)$$

$$(b) \quad E^\pi(\beta_{\tau_N}^N(\pi)) = \sup_{t \in C^N(\pi)} E^\pi(\beta_t^N(\pi))$$

Proposition 2 (3.1), (3.2) で定義された $\{\beta_n^N(\pi)\}$, τ_N に対し

$$E^\pi(x_{\tau_N}) = E^\pi[\beta_{\tau_N}^N(\pi)] = \beta_1^N(\pi)$$

(注意) Proposition 1 および Proposition 2 は、この報告の本末の目的のために全く準備的で事がらに過ぎないが、これ自体また別の意味で非常に興味ある結果である。これらの証明はかなりの補助定理を必要とするのでここでは省略するが、 $\{\beta_n^N(\pi)\}$ -process が super Martingale 产生をもつことから Martingale system theory を用いることを付記しておく。

$$X^N \equiv ASAS \dots S \quad (2N \text{ factors})$$

Σ^* ; $P(X^N)$ の上の σ -field で各 Borel subset $B \subset X^N$ に対し、写像 $\nu(B)$; $P(X^N) \rightarrow R^1$ が Σ^* -可測となる最小の σ -field

Σ^* は、weak-topology に関する $P(X^N)$ の Borel field と一致する。これが分かる。

$$\Gamma^N \equiv \{(s, v) \mid s \in S, v = e_\pi(s) \text{ for some } \pi \in \Pi^N\}$$

Lemma 3.1 Γ^N は $S \cdot P(X^N)$ の Borel subset である。

次に $\{\beta_n^N(\pi)\}$ -process は類似な方法で $\{\beta_n^N(v)\}$ -process を構成する。

$v \in P(X^N)$ を任意に与え、 v の分解を考へる。

$$v = v_1 v_2 \cdots v_{2N},$$

ただし

$$v_1 \in P(A)$$

$$v_2 \in Q(S|A)$$

$$v_{2n+1} \in Q(A|ASAS \cdots AS) \quad (\text{2n factors})$$

$$v_{2n} \in Q(S|ASAS \cdots A) \quad (\text{2n-1 factors}).$$

$\{\beta_m^N(v)\}$ を次式で backward に構成する。

$$\begin{cases} \beta_m^N(v) = \max [x_n, v_{2m-1} v_{2m} \beta_{m+1}^N(v)], & m=1, 2, \dots, N-1 \\ \beta_N^N(v) = x_N \end{cases} \quad (3.3)$$

ここで次の仮定をおく。

(A1) つきの (a), (b), (c) のうち少くとも一つが成立する:

$$(a) -\infty < \beta_1^N(v) < \infty \quad \text{for } \forall v \in P(X^N)$$

$$(b) r \geq 0, \quad q \geq 0$$

$$(c) r \leq 0, \quad q \leq 0.$$

(注意) (A1) の (b), (c) はそれぞれ positive case, negative case である。また (a) は §5, §6 で扱う discounted case ではない。discounted case は今后, D-case と呼ぶ。

$$v^{N*} \equiv \sup_{\pi \in \Pi^N} \beta_1^N(\pi).$$

一般に v^{N*} は可測になるとは限らないが、次の Lemma が成立つ。

Lemma 3.2 (A1) を仮定すると、 v^{N*} は絶対可測である。

(略証)

$$v(A, v) = \beta_1^N(v)(A)$$

$$B_\lambda = P^N \cap \{(A, v) \mid v(A, v) > \lambda\}$$

とおくと $\sum^* \alpha_i v_i = 1$ から $v(A, v)$ が $SP(X^N)$ 上で (A, v) に関する可測なること、Lemma 3.1 により B_λ は $SP(X^N)$ の Borel subset にある。一方

$$v^{N*}(A) = \sup_{v \in P_A^N} v(A, v) \quad (\text{ただし } P_A^N \text{ は } P^N \text{ の } A\text{-section})$$

だから、 $C_\lambda = \{\alpha \mid v^{N*}(\alpha) > \lambda\}$ は B_λ の $S \mapsto$ projection である。従って C_λ は analytic set. 故に、絶対可測である。 (Q.E.D.)

Lemma 3.3 (A1) を仮定すると、任意の $p \in P(S)$ 、任意の $\epsilon > 0$ に対して

$$p\{\beta_1^N(\alpha) \geq v^{N*} - \epsilon\} = 1$$

あるような $\alpha \in \prod^N$ が存在する。

(略証)

Lemma 3.2 によると、任意の $p \in P(S)$ に対して Borel set $N_p \subset S$ と可測度数 v_0 が存在して

$$\mu(N_1) = 0 \iff v_0(s) = v^{N^*}(s) \text{ for } s \notin N_1.$$

つまり

$$\begin{aligned} \Gamma_\varepsilon &\equiv \Gamma^N \cap \left[\{(s, v) \mid s \notin N_1, v(s, v) > v_0(s) - \varepsilon\} \right. \\ &\quad \left. \cup \{(s, v) \mid s \in N_1\} \right] \end{aligned}$$

とおくと、 Σ^* の定義と Lemma 3.1 により、 Γ_ε は $SP(X^N)$ の Borel subset である。また

$$\Gamma_\varepsilon \text{ の } s\text{-section} = \emptyset \text{ for all } s \in S$$

は、 Γ^N と v^{N^*} の定義より明らか。故に G.W. Mackey の 可測陰関数の定理 ([5]の Theorem 6.3) により、Borel set $N_2 \subset S$ と可測関数

$$y = \gamma(s) : S \rightarrow P(X^N)$$

が存在して、 $\mu(N_2) = 0$ かつ $(s, \gamma(s)) \in \Gamma_\varepsilon$ for $s \notin N_2$.

Σ^* の定義と γ の可測性から、各 $B \subset X^N$ に対して

$\gamma(s)(B)$ は s に関する可測。故に $\mu(\cdot | s) = \gamma(s)(\cdot)$ とおけば、 $\mu \in Q(X^N | s)$ となる。

$$\mu = \mu_1 \mu_2 \cdots \mu_{2n},$$

$$\mu_1 \in Q(A | S)$$

$$\mu_{2n} \in Q(S | SA \cdots A) \quad (2n \text{ factors})$$

$$\mu_{2n+1} \in Q(A | SA \cdots S) \quad (2n+1 \text{ factors})$$

と今解出来る。かつ

$$\mu_{2n}(\cdot | s_1 a_1 \cdots s_n a_n) = g_n(\cdot | s_1 a_1 \cdots s_n a_n) \text{ for } s_i \notin N_2$$

となる。 $\hat{\pi} = \{\hat{\pi}_1, \hat{\pi}_2, \dots, \hat{\pi}_N\}$ を

$$\hat{\pi}_n = \begin{cases} \mu_{2n-1} & \text{for } s_1 \notin N_2 \\ \pi'_n & \text{for } s_1 \in N_2 \end{cases} \quad \text{ただし } \pi'_n \text{ は任意} \\ (n=1, 2, \dots, N)$$

で定義すれば、 $s_1 \notin N_1 \cup N_2$ のならば

$$U(s_1, \delta_{(A_1)}) = \beta_1^N(\hat{\pi})_{s_1} \geq U_0(s_1) - \varepsilon = U^{N*}(s_1) - \varepsilon.$$

しかも $p(N_1 \cup N_2) = 0$ だから

$$p\{\beta_1^N(\hat{\pi}) \geq U^{N*} - \varepsilon\} = 1 \quad (\text{Q.E.D.})$$

(注意) Lemma 3.3 の証明で Mackey の下僕度数の定理を用いたが、これが適用出来るためには $P(X^N)$ が位相空間として standard Borel space であれば十分である。
(S はもともと standard Borel space である) $P(X^N)$ がこの § のはじめに導入した weak-topology に関する standard Borel space であることは K. R. Parthasarathy によって [6] の第 2 章に詳しく述べられている。

つきの定理は 有限 stage の model に対する最適政策の存在定理であるが、実はこれがこの § で求めた (p, ε, δ) -optimal s -policy の存在定理のための一つの重要な base である。

Theorem 3.1 (A1) を仮定する。任意の $p \in P(S)$, 任意の $\varepsilon > 0$ に対して $(p, \varepsilon)^N$ -optimal t -policy が存在する。

(署証)

Lemma 3.3 により 任意の $p \in P(S)$, $\varepsilon > 0$ に対して

$$\mathbb{P} \left\{ \beta_1^N(\hat{\pi}) \geq v^{N*} - \varepsilon \right\} = 1$$

すなはち $\hat{\pi} \in \Pi^{N-1}$ がある。故に

$$\mathbb{P} \left\{ \beta_1^N(\hat{\pi}) \geq \beta_1^N(\pi) - \varepsilon \right\} = 1 \quad \text{for all } \pi \in \Pi^{N-1} \quad (3.4)$$

π を任意の policy として, $\tau_N(\hat{\pi})$, $\tau_N(\pi)$ をそれぞれ (3.2) で定義される $\hat{\pi}$ と π に対する stopping time とする。

Proposition 2 により, (3.4) は代入して

$$\mathbb{P} \left\{ E^{\hat{\pi}}(x_{\tau_N(\hat{\pi})}) \geq E^{\pi}(x_{\tau_N(\pi)}) - \varepsilon \right\} = 1 \quad \text{for } \pi \in \Pi^{N-1}$$

Proposition 1-(a) により

$$\mathbb{P} \left\{ E^{\hat{\pi}}(x_{\tau_N(\hat{\pi})}) \geq E^{\pi}(x_t) - \varepsilon \right\} = 1 \quad \text{for } (\pi, t) \in \Lambda^N.$$

(Q.E.D.)

$\varepsilon = \varepsilon'$ 仮定 (A1) は代入して, 仮定 (B1) も満たす。

(B1) つまり (a), (b), (c) のうち少なくとも一つが成立つこと。

(a) (A1) の (a) が“各 $N \geq 1$ について成立”

(b) $r \geq 0$, $g \geq 0$

(c) $r \leq 0$, $g \leq 0$

(注意) (B1) の (a) は D-case では常に満たされる。

Lemma 3.4 (B1) を仮定する。 ρ を $P(S)$ の任意の要素、 ε を任意の正数とする。各 $N \geq 1$ に対して $(\rho, \varepsilon)^N$ -optimal t -policy を $(\hat{\pi}^N, \hat{\gamma}_N)$ とする。(その存在は Theorem 3.1 により保証される)。このとき一般性を失うなく、 ρ 確率 1 で $\{E^{\hat{\pi}^N}(x_{\hat{\gamma}_N})\}$ は単調非減少であるとしてよい。

ここであらかじめ次の仮定をおく。

$$(A2) \quad \sup_{(\pi, t) \in \Lambda} E^\pi(x_t) < \infty$$

$$(A3) \quad r_k' = r_k'' - r_k''' \quad (k=1, 2, \dots)$$

ただし $r_k''' \geq 0$ で r_k', r_k'' は共に SAS 上の実数値 Baire 测度。

$$(A4) \quad E^\pi(x_t') < \infty \text{ for every } (\pi, t) \in \Lambda$$

$$\text{ただし } x_t' = \sum_{k=1}^{t-1} \alpha^{k-1} r_k' + \alpha^{t-1} g_t.$$

$$(A5) \quad \{(x_n')^-, n \geq 1\} \text{ が } \pi \text{ に関して一様可積分}$$

$$\text{ただし } x_n' = \sum_{k=1}^{n-1} \alpha^{k-1} r_k' + \alpha^{n-1} g_n.$$

(注意) (A2) は D-case では常にみたされる。また

(A3)において $r_k''' \geq 0$ ということは r_k''' が cost または cost を含む損失の部分を表わすと思えばよい。このとき r_k' には可測性の他に制限がないから、我々の

問題は cost だけを評価関数として cost を最小にする最適化問題をも含むことになる。この意味で仮定 (A3) はごく一般的な仮定である。

以上の仮定のもとで次の Lemma を得る。

Lemma 3.5 (B1), (A2) ~ (A5) を仮定する。各 N に
対する $(p, \varepsilon)^N$ -optimal t -policy を $(\hat{\pi}^N, \hat{t}_N)$ で表わせば

$$\left| \left\{ \lim_{N \rightarrow \infty} E^{\hat{\pi}^N}(x_{\hat{t}_N}) \geq E^\pi(x_t) - \varepsilon \right\} \right| = 1 \quad \text{for } \forall (\pi, t) \in \Lambda$$

が成立する。

(署証)

一般性を失うことなく、 $E^\pi(x_t) \neq -\infty$ とする。

(A4) より $-\infty < E^\pi(x_t) \leq E^\pi(x'_t) < \infty$

$$\therefore -\infty < E^\pi(x''_t) < \infty \quad \text{ただし } x''_t = \sum_{k=1}^{t-1} \alpha^{t-k-1} r_k \quad (3.5)$$

一方 $t \in C(\pi)$ に対して

$$t > n \Rightarrow x_m \geq -(x'_m)^- - x''_t$$

右辺に

$$\int_{\{t > n\}} x_m^- d\pi \leq \int_{\{t > n\}} [(x'_m)^- + x''_t] d\pi \quad (3.6)$$

(3.5) (3.6) (A5) によると

$$\lim_{n \rightarrow \infty} \int_{\{t > n\}} x_m^- d\pi = 0 \quad (3.7)$$

次に, $t_N = \min(t, N)$ とおけば, $(\pi, t_N) \in \Lambda^N$.

故に $(\hat{\pi}^N, \hat{\tau}_N)$ の定義により

$$\begin{aligned}
 \int_{\{t \leq N\}} x_t d\pi &= E^\pi(x_{t_N}) - \int_{\{t > N\}} x_N d\pi \\
 &\leq E^{\hat{\pi}^N}(x_{\hat{\tau}_N}) + \varepsilon - \int_{\{t > N\}} x_N d\pi \quad w.p.1 \\
 &\leq E^{\hat{\pi}^N}(x_{\hat{\tau}_N}) + \varepsilon + \int_{\{t > N\}} x_N^- d\pi \quad w.p.1
 \end{aligned}$$

(3.8)

しかも (3.8) は各 $N \geq 1$ に対して成立。 (3.8) において

$N \rightarrow +\infty$ とすれば (A2), Lemma 3.4, (3.7) が得られる。

$\lim_{N \rightarrow \infty} E^{\hat{\pi}^N}(x_{\hat{\tau}_N})$ は存在して finite で, Lemma 9 により

結果が得られる。 (Q.E.D.)

この章の主目的である (p, ε, δ) -optimal s -policy の存在定理としてつきが与えられる。

Theorem 3.2 (B1), (A2)~(A5) を仮定する。

各 $p \in P(S)$, 各 $\varepsilon > 0$, 各 $\delta > 0$ に対して (p, ε, δ) -optimal s -policy が存在し, かつそれは t -policy の class 9 中にある。

(略証)

Lemma 3.5 は Egoroff の定理を適用すれば直ちに得られる。

Corollary 3.1 Theorem 3.2 と同じ仮定をおく。

(a) S が有限集合ならば、任意の $p \in P(S)$ 、任意の $\varepsilon > 0$
 \Rightarrow $\exists (p, \varepsilon)$ -optimal t -policy が存在する。

(b) $p \in P(S)$ が finite support をもてば、任意の $\varepsilon > 0$
 \Rightarrow $\exists (p, \varepsilon)$ -optimal t -policy が存在する。

(署証)

(a) または (b) の仮定があれば、 p -確率 1 で
 $\{E^{\hat{\pi}^N}(x_{\hat{\pi}_N})\}$ が一様収束することから明らかである。

§4. Markov stopped-policy

この § を通じて system の運動法則は Markov type であるとする。即ち

$$q_n \in Q(S|SA) \quad \text{for each } n$$

とする。ただし $Q(S|SA)$ における $|$ の右側の SA は n 番 stage における状態空間と行動空間を表わすものとする。

§3 で (p, ε, δ) -optimal t -policy の存在を示したが、この § では更にそれが Markov t -policy の class の中にあることを示す。そのためには、次のように $\{\beta_n^N(\pi)\}$ を $\{v_n^N(\pi)\}$ に変換する。

$$\alpha^{n-1} v_n^N(\pi) = \beta_n^N(\pi) - \sum_{k=1}^{n-1} \alpha^{k-1} \gamma_k, \quad n=1, 2, \dots, N \quad (4.1)$$

従って $\{\beta_n^N(\pi)\}$ の定義式 (3.1) は $\{v_n^N(\pi)\}$ を用いて次の
ように書き代えられる。

$$\begin{cases} v_n^N(\pi) = \max [g_n, \pi_n g_n \{ \alpha v_{n+1}^N(\pi) + \gamma_n \}], & n=1, 2, \dots, N-1 \\ v_N^N(\pi) = g_N \end{cases} \quad (4.2)$$

(注意) (4.2) をみると $\{\beta_n^N(\pi)\}$ -process π は、はつきり
(なかつた事が明らかになる。) 且つ (4.2) は policy π に
従って行動してゆくとき、最適な stopping rule によって
得られる optimal return に関する dynamic programming
form の再帰式 π あり、 $v_n^N(\pi)$ は残りの stage までか
n であるときの optimal return を表す。

次にあける Lemma 17 は π を引用した Mackey の
関数の定理に匹敵する重要な π への可測関数の
定理である。

Lemma 4.1 (Blackwell and Ryll-Nardzewski [1])

X, Y を Borel set とする。任意の $f \in Q(Y|X)$, 任意
の $w \in M(XY)$, 任意の $\epsilon > 0$ に対して

$$f[\{y; w(x, f(x)) \geq w(x, y) - \epsilon\} | x] = 1 \quad \text{for } \forall x \in X$$

なる Baire 関数 $f: X \rightarrow Y$ が存在する。

次の Lemma もまた 一つの Baire 因数の定理である。

Lemma 4.2 任意の $\pi, \beta \in P(S)$, $w \in M(SA)$, $\varepsilon > 0$

および各 n ($1 \leq n \leq N-1$) に対して

$$\beta \pi_1 g_1 \cdots \pi_{n-1} g_{n-1} \{ f_n w \geq w - \varepsilon \} = 1$$

なる Baire 因数 $f_n : S \rightarrow A$ が存在する。

(略証)

$$M \equiv \beta e_\pi \in P(SASA\cdots).$$

以下 " a_n が" 与えられたとの条件のもとで" の a_n に
ついての conditional prob. measure を π_n^* とするとき、

Lemma 4.1 より 任意の $w \in M(SA)$, $\varepsilon > 0$ および各
 n に対して

$$\pi_n^* \{ w(s_n, f_n(s_n)) \geq w(s_n, a_n) - \varepsilon \mid s_n \} = 1$$

for $\forall s_n \in S$

なる Baire 因数 $f_n : S \rightarrow A$ が存在する。すなは

$$\begin{aligned} M \{ w(s_n, f_n(s_n)) \geq w(s_n, a_n) - \varepsilon \} \\ = \beta \pi_1 g_1 \cdots \pi_{n-1} g_{n-1} \pi_n^* \{ w(s_n, f_n(s_n)) \geq w(s_n, a_n) - \varepsilon \\ \mid s_n \} \\ = \beta \pi_1 g_1 \cdots \pi_{n-1} g_{n-1} \cdot 1 = 1 \end{aligned}$$

(Q.E.D.)

Theorem 4.1 任意の $\beta \in P(S)$, $\varepsilon > 0$ および任意の t -
policy (π, t) に対して, (π, t) を (β, ε) -dominate すば Markov

t -policy が存在する。

(署意)

$(\pi, t) \in$ 任意の t -policy とする \forall

$$(\pi, t) \in \Lambda^N \text{ for some } N.$$

(4.1) から

$$\beta_i^N(\pi) = v_i^N(\pi).$$

Proposition 1, 2 は

$$E^\pi(x_t) \leq E^\pi(x_{\tau_N(\pi)}) = \beta_i^N(\pi) = v_i^N(\pi) \quad (4.3)$$

(4.2) は

$$\begin{aligned} v_{N-1}^N(\pi_{N-1}) &= \max [g_{N-1}, \pi_{N-1} g_{N-1} + \alpha g_N + r_{N-1}] \\ &= \max [g_{N-1}, \pi_{N-1} w_{N-1}]. \end{aligned} \quad (4.4)$$

\vdash

$$w_{N-1} = g_{N-1} + \alpha g_N + r_{N-1} \in M(SA).$$

$$\zeta \equiv \varepsilon / (1 + \alpha + \alpha^2 + \dots + \alpha^{N-1}) \quad \text{を}.$$

Lemma 4.2 は \forall Baire 海数 $f_{N-1} : S \rightarrow A$ が

$$\pi_{N-1} w_{N-1} \leq f_{N-1} w_{N-1} + \zeta \quad \text{with } p \pi_1 g_1 \dots \pi_{N-2} g_{N-2} - \text{prob. 1}$$

左 \vdash (4.4) から

$$\begin{aligned} v_{N-1}^N(\pi_{N-1}) &\leq \max [g_{N-1}, f_{N-1} w_{N-1} + \zeta] \\ &\leq \max [g_{N-1}, f_{N-1} w_{N-1}] + \zeta \quad " " " \\ &= v_{N-1}^N(f_{N-1}) + \zeta \quad " " " \end{aligned}$$

$$\begin{aligned} & \because \pi_{N-2} f_{N-2} \left\{ \alpha v_{N-1}^N(\pi_{N-1}) + r_{N-2} \right\} \\ & \leq \pi_{N-2} w_{N-2} + \alpha \zeta \quad \text{with } p \pi_1 f_1 \dots \pi_{N-3} f_{N-3} \text{-prob. 1} \end{aligned}$$

$\tau = \tau^*$

$$w_{N-2} = f_{N-2} \left\{ \alpha v_{N-1}^N(f_{N-1}) + r_{N-2} \right\} \in M(SA).$$

再び Lemma 4.2 により Baire 测度 $f_{N-2}; S \rightarrow A$ が“ τ ”とす

$$\begin{aligned} v_{N-2}^N(\pi_{N-2}, \pi_{N-1}) &= \max \left[g_{N-2}, \pi_{N-2} f_{N-2} \left\{ \alpha v_{N-1}^N(\pi_{N-1}) + r_{N-2} \right\} \right] \\ &\leq \max \left[g_{N-2}, f_{N-2} w_{N-2} + \zeta(1+\alpha) \right] \quad \text{with } p \pi_1 f_1 \dots \pi_{N-3} f_{N-3} \text{-prob. 1} \\ &\leq v_{N-2}^N(f_{N-2}, f_{N-1}) + \zeta(1+\alpha) \quad " " \end{aligned}$$

以下同様に Baire 测度 f_{N-3}, \dots, f_1 をえらんで

$$v_1^N(\pi) \leq v_1^N(f_1, f_2, \dots, f_{N-1}) + \zeta \quad \text{with } p \text{-prob. 1} \quad (4.5)$$

$$\pi^* = \{f_1, f_2, \dots, f_{N-1}\}$$

$$\tau^* = \text{the first } n \text{ such that } v_n^N(\pi^*) = g_n$$

$\times \# \times$

$$\tau^* = \tau_N(\pi^*). \quad (\text{たなびく } \tau_N(\pi^*) \text{ は (3.2) で定義され})$$

かつ、 (π^*, τ^*) は Markov t-policy である。

(4.3) により (4.5) を書き直して

$$p \left\{ E^{\pi^*}(x_{\tau^*}) \geq v_1^N(\pi) - \zeta \right\} = 1$$

再び (4.3) により

$$p \left\{ E^{\pi^*}(x_{\tau^*}) \geq E^\pi(x_\tau) - \zeta \right\} = 1.$$

(Q.E.D.)

Theorem 4.2 Theorem 3.2 と同じ仮定をおく。

任意の $p \in P(S)$, 任意の $\varepsilon > 0$, 任意の $\delta > 0$ に対して,
 (p, ε, δ) -optimal Markov t -policy が存在する。

(略証)

Theorem 3.2 と Theorem 4.1 より明らか。

Corollary 4.1 Theorem 3.2 と同じ仮定をおく。

- (a) S が有限集合であれば, 任意の $p \in P(S)$, $\varepsilon > 0$. 1=対して, (p, ε) -optimal Markov t -policy が存在する。
- (b) $p \in P(S)$ の finite support をもてば, 任意の $\varepsilon > 0$ 1=対して, (p, ε) -optimal Markov t -policy が存在する。

(注意) Theorem 4.1 の證明中で, $v_m^N(\pi)$ と書くべき所を
 $v_{N-1}^N(\pi_{N-1})$, $\pi_{N-2}^N(\pi_{N-2}, \pi_{N-1})$ 等と書いたが, 實は $v_m^N(\pi)$
 は $\pi = \{\pi_1, \pi_2, \dots\}$ の成るの中, $(\pi_m, \pi_{m+1}, \dots, \pi_{N-1})$
 のみに depend するからである。

§5 Stationary stopped-policy

この § では, D-case に対して, stationary s -policy の存在と, それに関連して再適方程式 (optimality equation)

等について論じる。

この § ではつきの仮定をおく。

- (i) $0 \leq \alpha < 1$
- (ii) $M(X)$ は X 上のすべての有界 Baire 関数からなる空間を表す。 $M(X)$ におけるノルムを $\|u\| = \sup_{x \in X} |u(x)|$ で定義する。
- (iii) $g_1 = g_2 = g_3 = \dots = g \in Q(S|SA)$.

次に degenerate $f \in Q(A|S)$ に対して operator T_f ; $M(S) \rightarrow M(S)$ を次式で定義する。更に A_f を定義する。

$$T_f u = f g (r + \alpha u),$$

$$A_f u = \max [g, T_f u].$$

Markov policy $\pi^N = \{f_1, f_2, \dots, f_N\}$ ($N < \infty$) に対して operator U_{π^N}, L_{π^N} ; $M(S) \rightarrow M(S)$ を次式で定義する。

$$U_{\pi^N} u = \max_{1 \leq n \leq N} T_{f_n} u,$$

$$L_{\pi^N} u = \max [g, U_{\pi^N} u].$$

これらの operator に関する直ちに次の 2 つの Lemma を得る。

Lemma 5.1

- (a) $T_f, A_f, U_{\pi^N}, L_{\pi^N}$ は monotone operator である。

(b) $u \in M(S)$, 定数 c に対して

$$T_f(u+c) = T_f u + \alpha c, \quad U_{\pi^N}(u+c) = U_{\pi^N} u + \alpha c.$$

(c) $u \in M(S)$, 定数 $c > 0$ に対して

$$A_f(u+c) \leq A_f u + \alpha c, \quad L_{\pi^N}(u+c) \leq L_{\pi^N} u + \alpha c.$$

(d) $u \in M(S)$, 定数 $c < 0$ に対して

$$A_f(u+c) \geq A_f u + \alpha c, \quad L_{\pi^N}(u+c) \geq L_{\pi^N} u + \alpha c.$$

Lemma 5.2 (a) $T_f, A_f, U_{\pi^N}, L_{\pi^N}$ は $M(S)$ 上の contraction mapping で、 α が contraction coefficient である。

(b) $u, v \in M(S)$, Markov policy $\pi = \{f_1, f_2, \dots\}$ に対して

$$\lim_{n \rightarrow \infty} \|A_{f_1} A_{f_2} \dots A_{f_n} u - A_{f_1} A_{f_2} \dots A_{f_n} v\| = 0.$$

次に π^N -generated の定義を与える。

Markov policy $\pi^N = \{f_1, f_2, \dots, f_N\}$ に対して

$f(S \rightarrow A)$ が π^N -generated ; S の Borel トポジション S_1, S_2, \dots

S_N が有って $f = f_n$ on S_n

Markov policy $\hat{\pi}^N = \{g_1, g_2, \dots, g_N\}$ が π^N -generated ;

各 g_n が π^N -generated

$F(\pi^N)$; すべての π^N -generated function の集合

$G^N(\pi^N)$; Π^N に属するすべての π^N -generated Markov policy の集合

Lemma 5.3

- (a) $T_f u \leq U_{\pi^N} u$ for any $u \in M(S)$, for any $f \in F(\pi^N)$
 (b) $A_f u \leq L_{\pi^N} u$ for any $u \in M(S)$, for any $f \in F(\pi^N)$

Lemma 5.4 任意の $u \in M(S)$ に対して次の f と $f \in F(\pi^N)$ が存在する:

$$T_f u = U_{\pi^N} u \quad \text{and} \quad A_f u = L_{\pi^N} u.$$

Lemma 5.5 u_N^* を L_{π^N} の fixed point とする。

(a) 任意の $\varepsilon > 0$ に対して N を十分大きくとれば

$$\beta_1^{N+1}(\hat{\pi}^N) - \varepsilon \leq u_N^* \text{ for } \hat{\pi}^N \in G^N(\pi^N).$$

(b) 任意の $\varepsilon > 0$ に対して N を十分大きくとれば

$$\beta_1^{N+1}(f^{(N)}) \geq u_N^* - \varepsilon$$

を $f \in F(\pi^N)$ をえらぶことが出来る。たゞ $f^{(N)} = \{f, f, \dots, f\}$

(N factors)

(累証)

(a) $\hat{\pi}^N = \{\hat{f}_1, \hat{f}_2, \dots, \hat{f}_N\} \in G^N(\pi^N)$ とする。

Lemma 5.3 (b) より

$$A_{\hat{f}_i} u \leq L_{\pi^N} u \text{ for } 1 \leq i \leq N, \text{ for } u \in M(S)$$

$$\therefore A_{\hat{f}_N} u_N^* \leq L_{\pi^N} u_N^* = u_N^*$$

$$\therefore A_{\hat{f}_1} A_{\hat{f}_2} \cdots A_{\hat{f}_N} u_N^* \leq u_N^*$$

(5.1)

Lemma 5.2 (b) より N を十分大きく取れば

$$\begin{aligned} u_N^* &\geq A_{f_1}^* A_{f_2}^* \cdots A_{f_N}^* g - \varepsilon \\ &= v_1^{N+1}(\hat{\pi}^N) - \varepsilon = \beta_1^{N+1}(\hat{\pi}^N) - \varepsilon. \end{aligned}$$

(b) Lemma 5.4 (b) より $f \in F(\pi^N)$ があって

$$A_f u_N^* = L_{\pi^N} u_N^* = u_N^*.$$

$$\therefore (A_f)^N u_N^* = u_N^*$$

Lemma 5.2 (b) より N を十分大きく取れば

$$(A_f)^N g + \varepsilon \geq (A_f)^N u_N^* = u_N^* \quad (5.2)$$

したがって $(A_f)^N g = v_1^{N+1}(f(N)) = \beta_1^{N+1}(f(N))$ だから (5.2) より

$$\beta_1^{N+1}(f(N)) \geq u_N^* - \varepsilon.$$

(Q.E.D.)

Theorem 5.1 (A3), (A4), (A5) を仮定する。このとき
任意の $p \in P(S)$, 任意の $\varepsilon > 0$, 任意の $\delta > 0$ に対し (p, ε, δ)
-optimal stationary s -policy が t -policy として
存在する。

(略証)

$(\pi^N, \tau_{N_1}) \in \Lambda^N$ を $(p, \frac{\varepsilon}{3}, \delta)$ -optimal Markov
 t -policy とする。(存在は Theorem 4.2.)

§3 における議論から (π^N, τ_{N_1}) は $(p, \frac{\varepsilon}{3})^{N_1}$ -
optimal として得られる。

一方, 各 N に対する $(p, \frac{\varepsilon}{3})^N$ -optimal t -policy を

(π^N, τ_N) とすれば Lemma 3.5 より

$$\Pr\left\{E^{\pi^N}(x_{\tau_N}) \geq E^\pi(x_t) - \frac{\varepsilon}{3}\right\} \geq 1 - \delta \quad \text{for } N \geq N_1,$$

for $(\pi, t) \in \Lambda \quad (5.3)$

\underline{u}_N^* の fixed point は u_{N-1}^* とすれば Lemma 5.5 (b) より

N_2 を十分大きくすれば、各 $N \geq N_2$ に対して

$$\exists f \in F(\pi^N); \beta_1^N(f^{(N-1)}) \geq u_{N-1}^* - \frac{\varepsilon}{3} \quad (5.4)$$

Lemma 5.5 (a) より N_3 を十分大きくすれば、各 $N \geq N_3$ に対して

$$u_{N-1}^* \geq \beta_1^N(\pi^N) - \frac{\varepsilon}{3} \quad (5.5)$$

$N_0 = \max [N_1, N_2, N_3]$ とおけば (5.4) (5.5) より

$$\beta_1^{N_0}(f^{(N_0-1)}) \geq u_{N_0-1}^* - \frac{\varepsilon}{3} \geq \beta_1^{N_0}(\pi^{N_0}) - \frac{2}{3}\varepsilon \quad (5.6)$$

$\tilde{\tau}_{N_0}$

$\tilde{\tau}_{N_0} = \text{the first } n \text{ such that } \beta_m^{N_0}(f^{(N_0-1)}) = x_m$

とおけば

$$\tilde{\tau}_{N_0} \in C^{N_0}(f^{(N_0-1)}) \Rightarrow \beta_1^{N_0}(f^{(N_0-1)}) = E^{f^{(N_0-1)}}(x_{\tilde{\tau}_{N_0}})$$

より (5.6) より

$$E^{f^{(N_0-1)}}(x_{\tilde{\tau}_{N_0}}) \geq E^{\pi^{N_0}}(x_{\tau_{N_0}}) - \frac{2}{3}\varepsilon \quad (5.7)$$

(5.3) (5.7) より

$$\Pr\left\{E^{f^{(N_0-1)}}(x_{\tilde{\tau}_{N_0}}) \geq E^\pi(x_t) - \varepsilon\right\} \geq 1 - \delta \quad \text{for } (\pi, t) \in \Lambda$$

$E^{f^{(N_0-1)}}(x_{\tilde{\tau}_{N_0}})$ は $(\rho, \varepsilon, \delta)$ -optimal stationary t -policy

である。 (Q.E.D.)

Corollary 5.1 Theorem 5.1 の同じ仮定をもく。

- (a) S が有限集合なら、任意の $p \in P(S)$, $\varepsilon > 0$ に対して
 (p, ε) -optimal stationary t-policy が“存在する。
- (b) $p \in P(S)$ が“finite support ではない”，任意の $\varepsilon > 0$
 に対して、 (p, ε) -optimal stationary t-policy が“存在する。

次の Markov policy $\pi = \{f_1, f_2, \dots\}$ に対して operator
 U_π, L_π を定義する。BTS

$$U_\pi u = \sup_n T_{f_n} u, \quad L_\pi u = \max[g, U_\pi u].$$

このとき、 U_π, L_π は monotone operator で、 π が $M(S)$ 上の
 contraction mapping である。

Markov policy $\pi = \{f_1, f_2, \dots\}$ に対して

$f : S \rightarrow A$ が π -generated ; S の Borel partition S_1, S_2, \dots
 “有る τ で $f = f_n$ on S_n ”

Markov policy $\hat{\pi} = \{g_1, g_2, \dots\}$ が π -generated ;

各 g_n が π -generated

$F(\pi)$; すべての π -generated function の集合

$G(\pi)$; π に属するすべての π -generated M-policy の集合

$$\Lambda_\pi = \{(\hat{\pi}, \hat{\tau}) \mid \hat{\pi} \in G(\pi), \hat{\tau} \in C(\hat{\pi})\}$$

Optimality equation

$A_a \pi f = a \Rightarrow f$ は π に対する operator である。

$\forall u \in M(S)$ に対して

$$u = \sup_{a \in A} A_a u \quad (5.8)$$

が成り立つ。 u は optimality equation を満たす π 。

(5.8) は optimality equation である。

Optimal return

$u \in M(S)$ が

$$(i) \quad u \geq E^{\hat{\pi}}(x_{\hat{t}}) \quad \text{for } \forall (\hat{\pi}, \hat{t}) \in \Lambda_{\pi}$$

$$(ii) \quad \forall \varepsilon > 0, \exists (\hat{\pi}, \hat{t}) \in \Lambda_{\pi}; E^{\hat{\pi}}(x_{\hat{t}}) \geq u - \varepsilon$$

を満たすとき、 u は optimal return in Λ_{π} である。

$u \in M(S)$ が

$$(iii) \quad u \geq E^{\hat{\pi}}(x_{\hat{t}}) \quad \text{for } \forall (\hat{\pi}, \hat{t}) \in \Lambda$$

$$(iv) \quad \forall \varepsilon > 0, \exists (\hat{\pi}, \hat{t}) \in \Lambda; E^{\hat{\pi}}(x_{\hat{t}}) \geq u - \varepsilon$$

を満たすとき、 u は optimal return である。

Theorem 5.2 (A3), (A4), (A5) をおく。

(a) π を任意の policy, Λ_{π} の fixed point は u^* である。
 u^* は optimal return in Λ_{π} である。

(b) $\varepsilon \geq 0$ とする。 ε -optimal s -policy がある。

各 $\varepsilon' > 0$ に対して $(\varepsilon/(1-\alpha) + \varepsilon')$ -optimal stationary s -policy π^* 存在する。

(c) $u \in M(S)$ が " $Aa u \leq u$ for all $a \in A$ をみたす" $E^\pi(x_t) \leq u$ for $(\pi, t) \in \Lambda$ が成立する。

(d) 各 $\varepsilon > 0$ に対して ε -optimal s -policy π^* が存在し、
optimal return は optimality equation をみたす。

(証明は非常には長いので省略。)

π が "semi-Markov" ; $\pi = \{\pi_1, \pi_2, \dots\}_{t=1}^\infty$ で $\pi_m \in Q(A|S)$
かつ π_m は degenerate

semi-Markov policy を用いて、次の定理が導かれる。

Theorem 5.3 (A3) (A4) (A5) をおく、

s -policy (π^*, t^*) が "optimal"

$\Leftrightarrow E^{\pi^*}(x_{t^*})$ が optimality equation をみたす

Theorem 5.4 optimality equation は $\mathbb{R}^n \rightarrow \mathbb{R}$ bounded
solution をもつ。

§ 6 Additional results.

ただし、action space が "countable set" または "finite set"
のときを扱う。

Theorem 6.1 (A3) (A4) (A5) をおく。

(a) A が countable set であれば、任意の $\epsilon > 0$ に対して ϵ -optimal stationary t -policy が存在する。

(b) (A3) (A5) の他に ((A4) は不要), (i) $\lim_{n \rightarrow \infty} x_n'' = \lim_{n \rightarrow \infty} \sum_{k=1}^n \alpha^{k-1} r_k'' = +\infty$, (ii) $\{x_n\}$ が各 e_n について一木尔可積分, (iii) $\sup_N E^{f^\infty}(x'_{T_N f^\infty}) < \infty$ for f^∞ の仮定をおく。このとき A が finite set なら stationary stopping time をもつより optimal stationary s -policy が存在する。

[付記]

この報告は [2] からの抜粋で、ここに記載しなかったかなりの部分、および証明を除いて箇所については [2] を見て頂きたい。

直前 §6 は optimal 又は ϵ -optimal s -policy の存在定理であるが、action space が countable ではなくても、一般に compact であれば、他に transition probability π に弱連続の仮定を加えれば Theorem 6.1 の結果と同様の結果が得られることが、最近 S. Iwamoto によって証明された。([3])

REFERENCES

- [1] D. Blackwell and C. Ryll-Nardzewski ; Non-existence of everywhere proper conditional distributions. Ann. Math. Statist. 34 (1963), 223-225.
- [2] N. Furukawa and S. Iwamoto ; Stopped decision processes on complete separable metric spaces.
(to appear)
- [3] S. Iwamoto ; Stopped decision processes on compact metric spaces. (to appear)
- [4] H. J. Kushner ; Computational procedures for optimal stopping problems for Markov chains. Jour. Math. Anal. Appl. 25 (1969), 607-615.
- [5] G. W. Mackey ; Borel structure in groups and their duals. Trans. Amer. Math. Soc. 85 (1957), 134 - 165.
- [6] K. R. Parthasarathy ; Probability measures on metric spaces. (1967), Academic Press.