

精密割引系をもつマルコフ決定過程について

和歌山大教育 門田良信

§ 1. 精密割引系について

無限段階のマルコフ決定過程において最もよく研究されている最適系に、 β -割引系と平均割引系がある。これらの系では β -割引期待利得 ($0 \leq \beta < 1$) が有限段階のその極限によって表わされるのに対して、平均期待利得は、最初の有限段階の期待利得とは無関係に、推移確率の時間に関する極限分布に依存して決まる値である。従って多くの場合、 β -割引最適政策の存在は平均最適政策のそれよりも緩い条件から導かれる事になる。更に、前者の存在定理を使って後者の存在を示す事も可能である。Blackwell [3] は、状態空間と決定空間が有限な場合には、1 に十分近いすべての β に対して β -割引最適定常政策が存在する事を示した。この政策による時刻に関する利得の総和は、他の政策によるそれよりも小さくはない事から、この政策が平均最適となる事は勿論、それよりも厳しい条件を持つ系に属するものである事がわかる。

Miller-Weinott [5] および Weinott [6] はこの政策を ∞ -割引最適政策と呼び、存在定理の別証明とそれを求めるための α

ルゴリズムを与えている。彼等の方法の特徴は、任意の定常政策に対応した利得を持つマルコフ連鎖を Laurent 展開し、その係数を辞書式順序に従って大小を比較する事により、定常政策の優劣を決める事にある。初項を -1 として第 n 項目までの係数比較において最大な政策を n -割引最適政策 ($n = -1, 0, 1, \dots$) と呼び、その n に関する極限をもって ∞ -割引最適政策を得たのである。これらの系を総称して *sensitive discount optimality criteria* と言う。(この報告では仮りに精密割引系と呼ぶ。)

但し、状態や決定の空間が有限でない場合には、事情は複雑になってくる。各系における最適政策は存在しない事もあるし、Jilynn [10] によれば ∞ -割引最適であっても平均最適とは限らない。精密割引系の最適政策の存在を示すすべての文献に共通した方法は、何らかの意味でエルゴード条件を設定する事によって Laurent 展開を導き、これを利用する事である。状態空間が可算な場合にはこのようにして、Hordijk [8], Hordijk-Sladky [11], Kadota [14, 15] によってそれぞれ、 0 -割引最適、 n -割引最適、 0 および ∞ -割引最適政策の存在が導かれている。状態空間が一般の場合にも、Taylor [9] は非周期的マルコフ連鎖の Laurent 展開を与え、Wijngaard [13] は 0 -割引最適定常政策の存在を導いている。また Sheu-Farn [16] は

決定空間をコンパクト空間として、0-割引最適定常政策の存在を示している。これらの存在定理の相異は、先に述べたエルゴード条件と、Laurent展開の係数の政策に関する連続性を保障する条件とに基づいていると思われる。その意味では[8], [11], [13]が与えた推移確率に関する条件は、[5], [6]の自然な拡張にはなっていない。

以下においては、精密割引最適政策に関する1つの存在定理を与える事を目的とする。取り扱うマルコフ決定過程の状態空間は、第2可算公理を満たす完備距離空間のボレル部分集合であり、決定空間は有限集合とする。諸定義と基礎的概念を§2で与えた後に、§3ではLaurent展開を導く。その方法は基本的には[5]と同じである。§4では ∞ -割引最適定常政策の存在を導く。また平均最適政策は-1割引最適政策に一致する事になる。証明の方法は、 β -割引最適政策の $\beta \rightarrow 1$ とした極限をLaurent展開を使って考える事による。§3, 4で設定される条件は[14]によって与えられたものである。これは有限状態空間では常に満たされているから、存在定理は[5]のそれを特別な場合として含む事になる。

§ 2. 準備

初めに記号の説明をする。可測空間を (X, \mathcal{F}) とする。 X が距離空間になっていれば、 \mathcal{F} は X のボレル集合族 $\mathcal{B}(X)$ にとる。 \mathcal{F} 可測な実数値有界関数の全体を $\mathcal{B}(X)$ と表わす。任意の $f \in \mathcal{B}(X)$ のノルムを $\|f\| = \sup\{|f(x)|; x \in X\}$ によって与える。2つの可測空間 $(X, \mathcal{F}), (Y, \mathcal{G})$ が与えられた時、 $\mathcal{F} \times \mathcal{G}$ 上の実数値有界関数で、任意の $y \in Y$ を固定すれば \mathcal{F} 上の \mathcal{G} -加法的集合関数、任意の $E \in \mathcal{F}$ を固定すれば \mathcal{G} 可測関数となるものの全体を $\mathcal{B}(X|Y)$ と表わす。任意の $H \in \mathcal{B}(X|Y)$ のノルムを $\|H\| = \sup\{|H(E|y)|; E \in \mathcal{F}, y \in Y\}$ によって与える。特に $Y = X, \mathcal{F} = \mathcal{G}$ のとき $P \in \mathcal{B}(X|X)$ が、任意の $x \in X, E \in \mathcal{F}$ に対して $P(E|x) \geq 0$ か $P(X|x) = 1$ を満たせば、 P を推移確率関数と呼ぶ。推移確率関数の全体を $\mathcal{P}(X|X)$ と表わす。また、 X が距離空間の時、 2^X は X の空でない閉部分集合の族とする。いま可測空間 (S, \mathcal{F}) と距離空間 A が与えられた時、その直積空間 $\mathcal{H}_1 = SA, \mathcal{H}_n = SAS \dots AS$ ($2n-1$ 個) が得られる。上記の表わし方により、空間 $\mathcal{B}(A), \mathcal{B}(2^A), \mathcal{B}(S|S), \mathcal{B}(S|SA), \mathcal{P}(S|S), \mathcal{P}(S|SA), \mathcal{B}(S), \mathcal{B}(SA), \mathcal{P}(A|\mathcal{H}_n)$ 等が定義される。

また $H, G \in \mathcal{B}(X|Y)$ が、すべての $E \in \mathcal{F}, y \in Y$ について $H(E|y) = G(E|y)$ を満たせば、これを簡単に $H = G$ と表わす事にする。同様の表記を $\mathcal{B}(X)$ の元についても用いる。不等号につ

いても等号の場合に準ずる。

任意の $H, G \in \mathcal{B}(S|S)$ に対して $HG(E|A)$ を

$$HG(E|A) = \int_S G(E|t) H(dt|A)$$

(Lebesgue-Stieltjes 積分) によって定義する。 $HG \in \mathcal{B}(S|S)$ である。特に $H, G \in \mathcal{P}(S|S)$ ならば $HG \in \mathcal{P}(S|S)$ である。同様に $H \in \mathcal{B}(S|S)$ と $r \in \mathcal{B}(S)$ に対して $Hr \in \mathcal{B}(S)$ が定義される。次の補助定理によって示される $\mathcal{B}(S|S)$, $\mathcal{B}(S)$ に関する性質は、以下のセクションにおいて断りなしによく使われる。

補助定理 2.1.

(a) $\mathcal{B}(S|S)$, $\mathcal{B}(S)$ は共に Banach 空間となり、 $\mathcal{P}(S|S)$ は $\mathcal{B}(S|S)$ の閉部分集合となる。

(b) $G, H, K \in \mathcal{B}(S|S)$ に対して、 $(GH)K = G(HK)$ かつ $H(G+K) = HG + HK$, $(G+K)H = GH + KH$ が成立する。

(c) $I(E|A)$ の値を $\lambda \in E$ ならば 1, $\lambda \notin E$ ならば 0 と定義すれば、 $I \in \mathcal{B}(S|S)$ であり、任意の $H \in \mathcal{B}(S|S)$ に対して $HI = IH = H$ が成立する。

(d) $H, G \in \mathcal{B}(S|S)$ に対して $\|HG\| \leq \|H\| \|G\|$.

(e) $H_n, G_n \in \mathcal{B}(S|S)$, $n=1, 2, \dots$ に対して、 $\|H_n - H\| \rightarrow 0$, $\|G_n - G\| \rightarrow 0$ as $n \rightarrow \infty$ ならば $\|H_n G_n - HG\| \rightarrow 0$ as $n \rightarrow \infty$

証明. (a), (d) は明らか。(e) は (d) を使って示される。

(b), (c) は、 K あるいは H が階段関数となる時に成立する事

を示した後にその極限をとればよい。□

任意の $H \in \mathbb{B}(S|S)$ に対して $H^1 = H$, $H^n = H \cdot H^{n-1} = H^{n-1} \cdot H$ とすれば、(b)により H^n は定義されて $H^n \in \mathbb{B}(S|S)$ となる。特に $P \in \mathbb{P}(S|S)$ ならば $P^n \in \mathbb{P}(S|S)$ であり、 $P^n(E|A)$ は状態 $A \in S$ から出発して n 時刻後に状態の集合 E にシステムが到達する確率を示す。

マルコフ決定過程を、順序づけられた次の因子の組 $\mathcal{M} = (S, A, P, r)$ によって定義する。 S は第2可算公理を満たす完備距離空間のボレル部分集合、その上の σ -加法族 $\mathbb{B}(S)$ は相対位相に関するボレル集合族とする。 A は有限集合 (離散距離空間と考える) とし、 $A(\cdot)$ を S から 2^A への集合値関数とする。 $P \in \mathbb{B}(S|SA)$ か $r \in \mathbb{B}(SA)$ とする。任意の $A \in S$ に対して $A(A)$ は状態 A における使用可能な決定の集合を示す。 $a \in A(A)$ が定めれば、利得 $r(A, a)$ が与えられシステムは推移確率 $P(\cdot | A, a)$ に従って次の段階 $A' \in S$ と進む。この決定を各段階において行い利得のある意味での和が最大となるような決定の列を捜す事が、マルコフ決定過程の問題である。

任意の $(A_1, a_1, A_2, \dots, A_n) \in \mathcal{H}_n$ に対して $\pi(A(A_n) | A_1, a_1, \dots, A_n) = 1$ を満たす $\pi \in \mathbb{P}(A | \mathcal{H}_n)$ の全体を $\mathbb{P}(\{A(A)\} | \mathcal{H}_n)$ と表わす。 $n = 1, 2, \dots$ に対して $\pi_n \in \mathbb{P}(\{A(A)\} | \mathcal{H}_n)$ を任意に取ってできる列 $\pi = (\pi_1, \pi_2, \dots)$ を政策と呼ぶ。各 π_n は時刻 n における決定を

それまでの履歴 $(\Delta_1, a_1, \Delta_2, \dots, \Delta_n)$ を参考にしながら確率的に定めるものである。すべての n , $(\Delta_1, a_1, \Delta_2, \dots, \Delta_n) \in \mathcal{H}_n$ に対して $\pi_n(\{a_n\} | \Delta_1, a_1, \Delta_2, \dots, \Delta_n) = 1$ となる $a_n \in A(\Delta_n)$ が $\Delta_1, a_1, \Delta_2, \dots, a_{n-1}$ に無関係に取れる時、 $f_n(\Delta_n) = a_n$ と定義する。 f_n は S から $A \wedge$ の関数で任意の $\Delta \in S$ に対して $f_n(\Delta) \in A(\Delta)$ を満たす。このような f_n の列 $\pi = (f_1, f_2, \dots)$ をマルコフ政策と呼ぶ。特に f_n が n に無関係になっていれば $\pi = f$ と表わし、これを定常政策と呼ぶ。定常政策の全体を F で表わす。

$A(\cdot)$ が $\mathcal{B}(S)$ と $\mathcal{B}(2^A)$ に関して可測、即ち任意の $B \in 2^A$ に対して $A^{-1}(B) \in \mathcal{B}(S)$ ならば、上記の f も可測となる。従って F は S から $A \wedge$ の $f(\Delta) \in A(\Delta)$ を満たす可測関数の全体に一致する。この時与えられた $P \in \mathcal{B}(S|SA)$, $r \in \mathcal{B}(SA)$ に対して、 $P(f)(E|\Delta) = P(E|\Delta, f(\Delta))$, $r(f)(\Delta) = r(\Delta, f(\Delta))$ と表わせれば、 $P(f) \in \mathcal{B}(S|S)$, $r(f) \in \mathcal{B}(S)$ となる。

$r \in \mathcal{B}(SA)$ と $\pi_n \in \mathcal{P}(\{A(\Delta)\} | \mathcal{H}_n)$ に対して $\pi_n r = \sum_{a \in A(\Delta_n)} r(\Delta_n, a) \pi_n(\{a\} | \Delta_1, a_1, \dots, \Delta_n)$ とすれば、 $\pi_n r \in \mathcal{B}(\mathcal{H}_n)$ となる。 $P \pi_n r = \int_S \pi_n r(\Delta_1, a_1, \dots, \Delta_n) P(d\Delta_n | \Delta_{n-1}, a_{n-1})$ とすれば、 $P \pi_n r \in \mathcal{B}(\mathcal{H}_{n-1} \wedge A)$ となる。 $\pi_{n-1} P \pi_n r = \sum_{a_{n-1} \in A(\Delta_{n-1})} P \pi_n r(\Delta_1, a_1, \dots, \Delta_{n-1}, a_{n-1}) \pi_{n-1}(\{a_{n-1}\} | \Delta_1, a_1, \dots, \Delta_{n-1})$ とすれば、 $\pi_{n-1} P \pi_n r \in \mathcal{B}(\mathcal{H}_{n-1})$ となる。このようにして $\pi_1 P \pi_2 P \dots P \pi_n r \in \mathcal{B}(S)$ を得る。 $\pi_1 P \pi_2 \dots P \pi_n r(\Delta)$ は、状態 $\Delta \in S$ から出発して政策 $\pi = (\pi_1, \pi_2, \dots)$ を使った時、時刻 n における期待利得を

示す。これらの定義の詳細は Blackwell [4], Furukawa [7] にある。任意の政策 $\pi = (\pi_1, \pi_2, \dots)$ に対して、 ρ -割引総期待利得および平均期待利得をそれぞれ、

$$V_\rho(\pi) = \sum_{n=0}^{\infty} \rho^{n+1} \pi_1 P \pi_2 P \dots \pi_{n-1} P \pi_n r, \quad \chi(\pi) = \liminf_n \frac{1}{n} \sum_{k=1}^n \pi_1 P \pi_2 P \dots P \pi_k r$$

によって定義する。但し $0 \leq \rho < 1$ かつ $\beta = \frac{1}{1+\rho}$ とし、 ρ -割引とは実際には β -割引の事である。 $V_\rho(\pi), \chi(\pi) \in \mathbb{B}(S)$ である。

政策 π^* が ρ -割引最適であるとは、任意の政策 π に対して $V_\rho(\pi^*) \geq V_\rho(\pi)$ が成立する事と言う。同様に平均最適であるとは $\chi(\pi^*) \geq \chi(\pi)$ となる事と言う。 n -割引最適であるとは

$$(1) \quad \liminf_{\rho \rightarrow 0} \rho^{-k} (V_\rho(\pi^*) - V_\rho(\pi)) \geq 0, \quad k = -1, 0, 1, \dots, n$$

が成立する事であり、 ∞ -割引最適であるとはすべての $n = -1, 0, 1, \dots$ に対して n -割引最適となる事である。従って ∞ -割引最適であるとは、任意の政策 π と $\Delta \in S$ に対して $0 < \rho < \rho_0$ ならば $V_\rho(\pi^*)(\Delta) \geq V_\rho(\pi)(\Delta)$ となる定数 $\rho_0 = \rho_0(\pi, \Delta) > 0$ が存在する事である。状態空間 S が有限ならば、上記の ρ_0 は π, Δ に無関係に存在する事が [3], [5], [6] によって知られている。また式 (1) の意味は Laurent 展開を得た段階で更に明らかとなる。

§ 3. Laurent 展開.

このセクションでは (S, \mathcal{F}) を可測空間とし、 $P \in \mathcal{P}(S|S)$,

$r \in B(S)$ を 1 つづつ取って固定しておく。P に関する次の仮定は以下ずっと用いられる。

仮定 I. 次の (2) を満たす正の数 M と $\mathcal{F} \times S$ から実数への関数 P^* が存在する。すべての $n=1, 2, \dots$, $\Delta \in S$, $E \in \mathcal{F}$ に対して

$$(2) \quad \left| \sum_{k=0}^n (P^k(E|\Delta) - P^*(E|\Delta)) \right| \leq M.$$

但し $P^0 = I$ とする。

(2) の両辺を $n+1$ で割って $n \rightarrow \infty$ にとれば P^* は極限分布 $\lim_{n \rightarrow \infty} \frac{1}{n+1} \sum_{k=0}^n P^k$ となり、その収束は Δ, E に関して一様である。また

$$(3) \quad P^* = PP^* = P^*P = (P^*)^2$$

が成立し、 $P^* \in P(S|S)$ である。

Doob [2] によると P が Doëblin 条件を満たしていれば、また Yosida-Kakutani [1] によればそれと同値な命題として、 $r \rightarrow Pr$ が $B(S)$ から $B(S)$ への線型作用素として quasi-strongly completely conti. ならば、仮定 I は満たされる。

$0 \leq \beta < 1$, $\beta = \frac{1}{1+\rho}$ とし、任意の $\Delta \in S$, $E \in \mathcal{F}$ に対して

$$H_\rho(E|\Delta) = \sum_{n=0}^{\infty} \beta^{n+1} (P^n(E|\Delta) - P^*(E|\Delta)) \quad \text{かつ} \quad h_\rho(\Delta) = (H_\rho r)(\Delta)$$

と定義すれば、 $H_\rho \in B(S|S)$, $h_\rho \in B(S)$ である。 $V_\rho = \sum_{n=0}^{\infty} \beta^{n+1} P^n$ の Laurent 展開を得るために、まず $\rho \rightarrow 0$ のとき $\|H_\rho - H\| \rightarrow 0$ となる $H \in B(S|S)$ の存在も言う。

補助定理 3.1. 仮定 I のもとですべての $\rho > 0$ に対して、

$$\|H_\rho\| \leq M, \quad \|h_\rho\| \leq M \cdot \|\rho\| \quad \text{が成立する。}$$

証明は [14] と同じである。

補助定理 3.2. 仮定 I のもとで H_p は

$$(I - \beta P) H_p = H_p (I - \beta P) = \beta (I - P^*), \quad P^* H_p = H_p P^* = 0$$

を満たす $B(SIS)$ の唯一の解である。

証明. (3) を使えば前半は明らか。 $Y \in B(SIS)$ が解とすれば $(I - \beta P)(Y - H_p) = 0$. $\beta^k P^k$ を左から掛けて $k=0, 1, \dots, n-1$ までの和を取れば, $(I - \beta^n P^n)(Y - H_p) = 0$ を得る。ノルムを取れば, $n \rightarrow \infty$ とした時, $\|Y - H_p\| = 0$ を得る。よって $Y = H_p$. \square

定理 3.3. 仮定 I のもとで, $\rho \rightarrow 0$ ならば $\|H_p - H\| \rightarrow 0$ なる $H \in B(SIS)$ が存在する。またこの時 $\|h_p - H_r\| \rightarrow 0$ が成立する。

証明. $B(SIS)$ は Banach 空間だから $\{P_k\}$ を $P_k \searrow 0$ as $k \rightarrow \infty$ に取った時, $\{H_{P_k}\}$ がコーシー列となる事を示せばよい。 $H = \lim_{k \rightarrow \infty} H_{P_k}$ が $\{P_k\}$ の取り方によらず一意に定まる事は, 次の補助定理 3.4 (a) によって示される。

任意の $\varepsilon > 0$ に対して $N > \frac{9M^2}{\varepsilon}$ なる N を取り, この N に対して $1 - \beta_0 < \varepsilon N / 2(N+2)^3$ を満たすように β_0 を取る。 $\beta_0 < \alpha < \beta < 1$ なる α, β を取り $\alpha = \frac{1}{1+\alpha}$, $\beta = \frac{1}{1+\beta}$ とする。補助定理 3.2 より $H_p = \beta(I - P^*) + \beta P H_p$. この式を ℓ 回反復して用いると,

$$H_p = \sum_{n=0}^{\ell} \beta^{n+1} (P^n - P^*) + \beta^{\ell+1} P^{\ell+1} H_p, \quad \ell = 0, 1, 2, \dots$$

$\ell = 0, 1, \dots, N$ までを加えて $N+1$ で割ると,

$$(4) \quad H_p = \frac{1}{N+1} \sum_{k=0}^N \sum_{n=0}^k \beta^{n+1} (P^n - P^*) + \frac{1}{N+1} \sum_{k=0}^N \beta^{k+1} P^{k+1} H_p$$

$$x_k = \sum_{j=0}^k (P^j - P^*) H_p \quad \text{とおくと補助定理 3.1 と仮定 I に}$$

より $\|x_k\| \leq M^2$. よって

$$\begin{aligned} \left\| \frac{1}{N+1} \sum_{k=0}^N \beta^{k+1} P^{k+1} H_p \right\| &= \left\| \frac{1}{N+1} \sum_{k=0}^N \beta^{k+1} (P^{k+1} - P^*) H_p \right\| \\ &\leq \frac{1}{N+1} \left\| \sum_{k=0}^N \beta^{k+1} (x_{k+1} - x_k) \right\| < \frac{3M^2}{N} < \frac{\varepsilon}{3} \end{aligned}$$

H_α についても (4) と同様の式が成立するから、結局

$$(5) \quad \|H_p - H_\alpha\| < \frac{1}{N+1} \sum_{k=0}^N \sum_{n=0}^k |\beta^{n+1} - \alpha^{n+1}| \cdot \|P^n - P^*\| + \frac{2}{3}\varepsilon$$

$\|P^n - P^*\| \leq 2$ および $\beta^{n+1} - \alpha^{n+1} < (1-\alpha) \left(\sum_{k=0}^n \alpha^k\right)$ を使って (5) の右辺第 1 項を評価してやれば

$$\|H_p - H_\alpha\| < (1-\alpha) \frac{2(N+2)^3}{3(N+1)} + \frac{2}{3}\varepsilon < \varepsilon$$

この事は $\{H_{p_k}\}$ がコーシー列となる事を示している。□

以下の補助定理 3.4、定理 3.5 に現われる H は定理 3.3 で求められたものである。3.4 (a) は補助定理 2.1 (e) および 3.2 を使って示される。他の証明も基本的には [5] と同じだから省略する。

補助定理 3.4. 仮定 I のもとで

(a) 補助定理 3.2 は H_p のかわりに H , $\beta=1$ として成立する。

(b) $\rho > 0$ に対して $M_\rho = \sum_{n=0}^{\infty} \rho^{n+1} P^n$ とおくと、 $M_\rho \in \mathbb{B}(S|S)$ で $P^* M_\rho = \frac{1}{\rho} P^*$, $M_\rho = P^* M_\rho + H_p$ が成立する。

(c) $\rho > 0$ に対して $L_\rho = \sum_{n=0}^{\infty} \rho^n (-1)^n H^n$ とおくと、 L_ρ は

$$(I + \rho H) L_\rho = L_\rho (I + \rho H) = I$$

も満たす $B(S)$ の唯一の解である。

(d) $\rho > 0$ に対して、 $H_\rho = L_\rho H = H L_\rho$ が成立する。

定理 3.5. 仮定 I のもとで

(a) $x = P^*r$ とおくと x は、 $(I - P)x = 0$ かつ $P^*x = x$ を満たす $B(S)$ の唯一の解である。また $y_0 = Hr$ とおくと y_0 は、 $(I - P)y = r - x$ かつ $P^*y = 0$ を満たす $B(S)$ の唯一の解である。

(b) $0 < \rho < \frac{1}{M}$ なる ρ に対して

$$(6) \quad V_\rho = \sum_{n=0}^{\infty} \rho^{n+1} P^n r = \sum_{n=-1}^{\infty} y_n \rho^n$$

が成立する。ここで $y_{-1} = x = P^*r$, $y_0 = Hr$, $n=1, 2, \dots$ に対しては $y_n = (-1)^n H^{n+1} r$ とする。 M は (2) を満たす定数とする。

系 3.6. 仮定 I のもとで

(a) $n = -1, 0, \dots$, $k = -1, 0, \dots, n$ に対して定理 3.5 (b) によって定義された $\{y_{-1}, y_0, \dots, y_n\}$ は

$$(7) \quad r_k = y_{k-1} + y_k - P y_k, \quad k = -1, 0, \dots, n$$

を満たす。但し、 $y_{-2} = 0$ かつ $k \neq 0$ ならば $r_k = 0$, $k = 0$ ならば $r_k = r$ とする。

(b) 逆に $n = 0, 1, \dots$ に対して $z_{-1}, z_0, \dots, z_n \in B(S)$ が (7) を満たせば、 $k = -1, 0, \dots, n-1$ に対して $y_k = z_k$ が成立する。

(定理 3.5 (a) は系 3.6 で $n = 0$ とした時に相当する。)

証明. (6) の左辺より $V_\rho = \beta r + \beta P V_\rho$ であり、 V_ρ に (6) の右辺を代入すれば、すべての $\rho > 0$ に対して、

$$0 = \beta \{ r + (P - (1+e)I) V_p \}$$

$$= \beta \left\{ \frac{1}{P} (P y_1 - y_1) + (r + P y_0 - y_0 - y_1) + \sum_{n=1}^{\infty} (P y_n - y_n - y_{n-1}) P^n \right\}$$

が成立する。よって P^n の各係数は 0 となり (a) が示された。
 逆に $k = n-1, n$ に対して (7) が成立して、 $y_{n-2} = z_{n-2}$ ならば、(7) の n に関する式に P^* を乗じて、 $P^* r_n = P^* y_{n-1} = P^* z_{n-1}$ 。また

$$0 = r_{n-1} - r_{n-1} = y_{n-1} - z_{n-1} - P(y_{n-1} - z_{n-1})$$

だから、 $y_{n-1} - z_{n-1} = P^*(y_{n-1} - z_{n-1}) = 0$ を得る。詳しくは n に関する帰納法によって示される。□

定理 3.5 の (a) は [3] の拡張であり、定理 3.5 (b) と系 3.6 は [5] の拡張である。[5], [6] により (6) の右辺は Laurent 展開と呼ばれている。 $\beta \rightarrow 1$ のとき $P \rightarrow 0$ だから (6) の右辺は V_p の漸近的性質を詳細に告げている。

§ 4. 存在定理

セクション 2 で定義したマルコフ決定過程 $\mathcal{M} = (S, A, P, r)$ に関して次の仮定をする。

仮定 II. $A(\cdot); S \rightarrow 2^A$ は $\mathcal{B}(S)$ と $\mathcal{B}(2^A)$ に関して可測。

仮定 III. すべての $f \in F$, $n = 1, 2, \dots$, に対して $\| \sum_{k=0}^n (P(f)^k - P(f)^*) \| \leq M$ をみたす定数 M が存在する。

仮定 III は、仮定 I で $P = P(f)$, $P^* = P(f)^*$ としたものに f

$\in F$ に関する一様有界性を付け加えたものである。 S が有限ならば仮定 II, III は消えてしまう。

$\{f_n\} \subset F$, $f \in F$ とし, 任意の $\Delta \in S$ に対して $\lim_{n \rightarrow \infty} f_n(\Delta) = f(\Delta)$ が成立すれば, $\{f_n\}$ は f に収束すると言う。決定空間 A は有限だからこの事は, $n \geq N$ ならば $f_n(\Delta) = f(\Delta)$ となる自然数 N が存在する事を意味する。

仮定 II, III より任意の $f \in F$ に対して $V_p(f) = \sum_{n=0}^{\infty} \rho^{n+1} P(f)^n r(f)$ の Laurent 展開が定理 3.5 (b) によって導かれる。この時の ρ^k の係数を $y_k = y_k(f)$ で表わす。次の補助定理は $\{y_k(f)\}$ の F 上での連続性を示すものである。

補助定理 4.1. 仮定 II, III のもとで $\{f_n\} \subset F$ が $f \in F$ に収束すれば, 任意の $\Delta \in S$ に対して

$$(8) \quad \lim_{n \rightarrow \infty} \sup_{E \in \mathcal{B}(S)} |P(f_n)^*(E|\Delta) - P(f)^*(E|\Delta)| = 0$$

かつ $k = -1, 0, 1, \dots$ に対して $\lim_{n \rightarrow \infty} y_k(f_n)(\Delta) = y_k(f)(\Delta)$ が成立する。

証明. $P(f)(\cdot|\Delta)$ と f_n にエゴロフの定理を適用する事により任意の $\varepsilon > 0$, $k = 2, 3, \dots$ に対して, $P(f)(S-K|\Delta) < \frac{\varepsilon}{k-1}$, $\Delta \in K$ かつ $t \in K$, $n \geq N$ ならば $f_n(t) = f(t)$ となる $K \in \mathcal{B}(S)$ と自然数 N が存在する。これより $n \geq N$ ならば

$$|P(f_n)^k(E|\Delta) - P(f)^k(E|\Delta)| \leq \sum_{j=1}^{k-1} P(f)^j(S-K|\Delta) < \varepsilon.$$

よって任意の $\Delta \in S$ に対して $P(f_n)^k$ は $P(f)^k$ に $E \in \mathcal{B}(S)$ に関して一様収束する。 $P(f)^*$ を仮定 III を使って $\frac{1}{k+1} \sum_{j=0}^k P(f)^j$ で評価する事

により (8) が示される。またこれらの事より $H_p(f_n)(E|A)$ が $n \rightarrow \infty$ とした時 $H_p(f)(E|A)$ に $E \in \mathcal{B}(S)$ に関して一様収束する事が言える。定理 3.3 により $p > 0$ を十分小さく取れば $\|H(f) - H_p(f)\| < \varepsilon$ だから、結局 $H(f_n)(E|A)$ も $H(f)(E|A)$ に $E \in \mathcal{B}(S)$ に関して一様収束する。これより $\lim_{n \rightarrow \infty} y_k(f_n)(A) = y_k(f)(A)$ は容易である。□

任意の $f \in F$, $A \in S$ に対して $Y_n(f)(A) = (y_{-1}(f)(A), y_0(f)(A), \dots, y_n(f)(A))$ と定める。 $f, g \in F$ についても $y_k(f)(A) \neq y_k(g)(A)$ となる $k = -1, 0, \dots, n$ があれば、そのような最小の k に対しては $y_k(f)(A) > y_k(g)(A)$ が成立する時、 $Y_n(f)(A) \succ Y_n(g)(A)$ と表わす。すべての $g \in F$, $A \in S$ に対して $Y_n(f)(A) \succ Y_n(g)(A)$ を満たす $f \in F$ の全体を D_n , $n = -1, 0, \dots$ によって表わす。任意の $A \in S$, $a \in A(A)$, $f \in F$ に対して、

$$\psi_n(A, a; f) = r_n(A, a) + P y_n(A, a) - y_n(f)(A) - y_{n-1}(f)(A)$$

$$\Psi_n(A, a; f) = (\psi_{-1}(A, a; f), \psi_0(A, a; f), \dots, \psi_n(A, a; f))$$

とする。 $f \in D_n$ のとき $A_n(A) = \{a \in A(A); \Psi_n(A, a; f) = 0\}$ とする。与えられた \mathcal{M} から各 $A \in S$ において使用可能な決定を $A_n(A)$ に制限したマルコフ決定過程を考え、これを \mathcal{M}_n と表わす。

系 3.6 (a) により $\Psi_n(A, f(A); f) = 0$ だから $A_n(A) \neq \emptyset$ である。また [7] の Lemma 4.5 により $A_n(\cdot)$ は $\mathcal{B}(S)$ と 2^A に関して可測となるから、 \mathcal{M}_n は仮定 II, III を満たす事になる。仮定 II のもとでは任意の $p > 0$ に対して p -割引最適定常政策が存在する事

が [7] によって知られている。A の元に適当に順番をつけ、 \mathcal{M}_n において $P_k \searrow 0$ に対して $F_n \ni f_k$ を P_k -割引最適定常政策とすれば、 $\liminf_{k \rightarrow \infty} f_k(A) = f(A)$ が存在する。 $f \in F_n$ となり、下極限政策と呼ぶ。ここで F_n は \mathcal{M}_n における定常政策の全体とする。

補助定理 4.2. 仮定 II, III のもとで $f \in \mathcal{D}_n$ ならば、 f は n -割引最適定常政策である。

証明. $n-1$ の時を仮定する。 $\liminf_P P^{-n} (V_P(f)(A) - V_P(\pi)(A)) = \alpha$ を与える $P_k \searrow 0$ をとり、 f_k を P_k -割引最適定常政策とする。 $Y_n(f) \geq Y_n(f_k)$ だから $y_l(f_k)(A) = y_l(f)(A)$ となる l が可算個あるような $l = -1, 0, \dots, n$ が存在すれば、そのような l の最大値と P_k, f_k を改めて l, P_k, f_k と表わす。

(9)
$$\liminf_{k \rightarrow \infty} P_k^{-l} (V_{P_k}(f)(A) - V_{P_k}(f_k)(A)) + \liminf_{k \rightarrow \infty} P_k^{-l} (V_{P_k}(f_k)(A) - V_{P_k}(\pi)(A))$$
 において第 2 項目は非負である。第 1 項目は、 $-1 \leq l \leq n-1$ ならば帰納法の仮定より非負、 $l=n$ ならば l のとり方により 0 である。もし $y_{-1}(f_k)(A) < y_{-1}(f)(A)$ ならば第 1 項目は $l=-1$ の時に正、 $l \geq 0$ の時、すべての y_n が有界である事により ∞ となる。いずれの場合にも $\alpha \geq 0$ だから、 n -割引最適である。□

補助定理 4.3. 仮定 II, III のもとで任意の $f \in \mathcal{D}_n$ を取って \mathcal{M}_n をつくと、 \mathcal{M}_n における下極限政策 f_* は \mathcal{D}_{n+1} に属する。但し $\mathcal{D}_{-2} = F, \mathcal{M}_{-2} = \mathcal{M}$ とする。

証明. 任意の $g \in F_n$ をとれば $U_n(A, g(A); f) = 0$ がすべての $A \in S$ について成立するから, 系3.6により $Y_{n-1}(g) = Y_{n-1}(f)$. よって補助定理4.2により g は $n-1$ -割引最適となる. 任意の $A \in S$ に対して $f_*(A) = \lim_{k \rightarrow \infty} f_k(A)$ となる $\{f_k\}$ を選ぶ. (9)において $l = n$, $f = f_*$ とした式を考えると, 第2項目は非負であり第1項目は (6) を使って $\lim_{k \rightarrow \infty} (Y_n(f_*)(A) - Y_n(f_k)(A)) = 0$ に等しい. よって $f_* \in D_n$ となる. $Y_n(f_*) \geq Y_n(f_k)$ を使えば同様にして $f_* \in D_{n+1}$ を得る. \square

定理4.4. 仮定II, IIIを満たすマルコフ決定過程 \mathcal{M} は ∞ -割引最適定常政策をもつ.

証明. 補助定理4.3により $D_n \neq \emptyset$ であり, 完備距離空間 F の閉集合となる. また $\{D_n\}$ は単調非増加な集合列であるから, $\bigcap_{n=2}^{\infty} D_n \neq \emptyset$ となる. 補助定理4.2により, この事は ∞ -割引最適定常政策の存在を示している. \square

参考文献

- [1] K. Yosida and S. Kakutani, "Operator-theoretical treatment of Markoff's process and mean ergodic theorem," *Ann. Math.* 42, No.1, 1941.
- [2] J. Doob, "Stochastic Processes," Wiley, New York, 1953.
- [3] D. Blackwell, "Discrete dynamic programming," *Ann. Math. Statist.* 33, 1962.

- [4] D. Blackwell, "Discounted dynamic programming," *Ann. Math. Statist.* 36. 1965.
- [5] B. L. Miller and A. F. Veinott, jr., "Discrete dynamic programming with a small interest rate," *Ann. Math. Statist.* 40. 1969.
- [6] A. F. Veinott, jr., "Discrete dynamic programming with sensitive discount optimality criteria," *Ann. Math. Statist.* 40. 1969.
- [7] N. Furukawa, "Markovian decision processes with compact action spaces," *Ann. Math. Statist.* 43. 1972.
- [8] A. Hordijk, "Regenerative Markov decision models," *Math. Programming Study* 6. 1976.
- [9] H. M. Taylor, "A Laurent series for the resolvent of a strongly continuous stochastic semi-group," *Math. Programming Study* 6. 1976.
- [10] J. Flynn, "Conditions for the equivalence of optimality criteria in dynamic programming," *Ann. Statist.* 4. 1976.
- [11] A. Hordijk and K. Sladky, "Sensitive optimality criteria in countable state dynamic programming," *Math. Oper. Res.* 2. 1977.
- [12] J. Wijngaard, "Stationary Markovian decision problems and perturbation theory of quasi-compact linear operators," *Math. Oper. Res.* 2. 1977.
- [13] J. Wijngaard, "Sensitive optimality in stationary Markovian decision problems on a general state space," *Math. Centre Tracts.* 93. 1977.
- [14] Y. Kadota, "Countable state Markovian decision processes under the

Doëblin conditions," *Bull. math. Statist.* 19, 1979.

[15] 門田. 「可算状態マルコフ決定過程の sensitive discount 最適系
について」 数理解析研究所講究録 358. 「決定過程論とその周辺」 1979.

[16] S. S. Shew and K. L. Farn, "A sufficient condition for the existence
of a stationary 1-optimal plan in compact action Markovian decision pro-
cesses," *Recent Developments in Markov Decision Processes*, Academic
Press, 1980.