

コンパクトな行動空間をもつ部分的観測  
セミマルコフ決定過程

長岡高専 涌田和芳 (Kazuyoshi Wakuta)

§ 1. 序

部分的観測マルコフ決定過程は多くの人により研究されてきた。その一般化として、White [5] は部分的観測セミマルコフ決定過程 (POSMDP) を導入した。ここでは、基礎過程として有限段階で状態が離散的に推移するセミマルコフ過程が用いられた。ここでは Wakuta [2,3] に引き続き、無限段階で状態が連続的に推移する場合を考えるが、状態空間は有限集合で、行動空間はコンパクト距離空間とする。

§ 2. POSMDP

まず記号について説明する。一般に、 $X$  と  $Y$  を空でないボレル集合とするとき、 $X$  上の確率測度の全体を  $P(X)$ 、 $X$  から  $Y$  への推移確率の全体を  $Q(Y|X)$  と表わす。また、任意の  $\mu \in P(X)$  と  $\nu \in Q(Y|X)$  に対して、 $\mu \nu$  は  $X \times Y$  で定義される通

常の確率測度とする。  $g$  が退化しているとき、すなわち、すべての  $x \in X$  に対して、  $g(\{x\} | x) = 1$  なる  $X$  から  $Y$  へのボレル可測写像  $u$  が存在するとき、  $g$  のかわりに  $\phi u$  とも書く。  
 $X$  上の有界ボレル可測関数の全体を  $M(X)$ 、有界連続関数の全体を  $C(X)$  と表わす。

POSM DP は次のもので定められる:  $S$  は有限集合で、システムの状態集合;  $M$  は有限集合で、可能な信号の集合;  $A$  はボレル集合で、行動集合;  $g \in Q(S \times R_+ | S \times A)$  はシステムの運動法則、ここで  $R_+ = [0, \infty)$ 、  $g$  は

$$(2.1) \quad g(\{j\} \times [0, t] | i, a) = \int_0^t k(j, \Delta | i, a) d\Delta,$$

$$i, j \in S, a \in A, t \geq 0$$

なる確率密度関数  $k$  をもつと仮定する。ただし、  $d\Delta$  はルバーク測度で、  $k(j, \cdot | i, \cdot)$ 、  $i, j \in S$ 、は  $(\Delta, a)$  の非負ボレル可測関数とする;  $r \in M(S \times A)$  は、利得関数;  $\alpha$  は正数で、割引因子。

有限時間区間で無限回の推移が起こらないように、次の仮定をおく。

Condition (1). (cf. Ross [1]) すべての  $(i, a) \in S \times A$  に対して、

$$(2.2) \quad g(S \times [0, \delta] | i, a) \leq 1 - \varepsilon$$

なる  $\varepsilon > 0$  と  $\delta > 0$  が存在する。

簡単のために、 $\Phi = P(S)$  とおく。すなわち、 $\Phi$  は  $S$  上のすべての確率測度の集合を表わす。 $\varphi \in \Phi$  の  $i$  番目の成分を  $\varphi(i)$  とすると、 $\varphi(i), i \in S$ , は  $\varphi$  について連続なので、 $\phi(i|\varphi) = \varphi(i), i \in S, \varphi \in \Phi$ , として  $\phi \in Q(S|\Phi)$  を定義することができる。行動を選択するための政策は、 $\omega = \{\omega_0, \omega_1, \dots\}$ ,  $\omega_n \in Q(A|H_n), n \geq 0$ , である。ここで、 $H_n = \Phi \times (A \times R_+ \times M)^n$  は  $n$  番目の行動が選択されるときの観測可能な歴史の集合であり、 $\varphi \in H_0 = \Phi$  はシステムの状態の初期情報を表わす。 $i_n, m_n, a_n, t_n$  は、それぞれ、 $n$  番目の状態、観測信号、行動、滞在時間を表わす。 $T_n = t_1 + \dots + t_n$  ( $T_0 = 0$ ) とおくと、 $T_n$  は、 $n$  番目の決定時刻を表わす。政策  $\omega$  と  $\phi, g, \gamma$  に対して、 $Q(S \times (A \times S \times R_+ \times M)^\infty | \Phi)$  の要素  $\phi_\omega = \phi \omega \cdot g \cdot \gamma \omega \cdot g \cdot \gamma \dots$  が定義できる。政策  $\omega$  を用いたときの期待合計割引利得を

$$(2.3) \quad J_\omega(\varphi) = E_\omega \left[ \sum_{n=0}^{\infty} e^{-\alpha T_n} r(i_n, a_n) | \varphi \right], \quad \varphi \in \Phi,$$

と定義する。ここで、 $E_\omega[\cdot | \varphi]$  は  $\phi_\omega\{\cdot | \varphi\}$  による条件付期待値を表わす。我々の最適化問題は、すべての政策の中で  $J_\omega$  を最大にすることである。すべての政策  $\omega$  とすべての  $\varphi \in \Phi$  に対して、 $J_{\omega^*}(\varphi) \geq J_\omega(\varphi)$  が成り立つとき、 $\omega^*$  は最適であるという。

Remark 1. 利得関数  $r$  は、次に推移する状態と滞在時間  
に依存するように一般化できる。

### § 3. 同値な SMDP の構成

POSMDP と同値な通常の (完全に観測可能な) セミマル  
コフ決定過程 (SMDP) を構成する。

政策を任意に固定し、 $h_n = (\varphi, a_0, t_1, m_1, \dots, a_{n-1}, t_n, m_n)$   
が与えられたときの  $i_n$  の条件付分布を  $\varphi_n$  とする。そして、 $\varphi_n$   
に対して、 $\varphi_n(i) = \varphi_n\{i_n = i \mid h_n\}$ ,  $i \in S$ , として  $\varphi_n \in \mathfrak{P}$  を対応  
させる。ベイスの公式より、 $\{\varphi_n\}_{n=0,1,2,\dots}$  の再帰式

$$(3.1) \quad \varphi_{n+1}(j) = \frac{\sum_{i \in S} k(j, t \mid i, a) g(m \mid j) \varphi_n(i)}{\sum_{i \in S} \sum_{j \in S} k(j, t \mid i, a) g(m \mid j) \varphi_n(i)}, \quad j \in S,$$

を得る。ここで、 $\varphi_0$  は初期情報である。そこで、写像  $u$ :  
 $D \rightarrow \mathfrak{P}$  を

$$(3.2) \quad u(\varphi, a, t, m)(j) = \frac{\sum_{i \in S} k(j, t \mid i, a) g(m \mid j) \varphi(i)}{\sum_{i \in S} \sum_{j \in S} k(j, t \mid i, a) g(m \mid j) \varphi(i)}, \quad j \in S,$$

として定義する。ここで、 $D = \{(\varphi, a, t, m) \mid \sum_{i \in S} \sum_{j \in S} k(j, t \mid i, a) g(m \mid j) \varphi(i) > 0\}$ 。 $u$  はボレル可測である。 $D^c$  は確率ゼ

口の集合なので、 $D^c$ 上に重の任意の値を与えて、 $u$ の定義域を全空間 $\Phi \times A \times R_+ \times M$ に拡張する。この拡張された写像 $\hat{u}$ は再びボレル可測である。 $\hat{u}$ を繰返し用いることにより、観測可能な履歴 $h_n = (\varphi, a_0, x_1, m_1, \dots, a_{n-1}, x_n, m_n)$ はボレル可測に $b_n = (\varphi, a_0, \varphi_1, x_1, \dots, a_{n-1}, \varphi_n, x_n) \in B_n = \Phi \times (A \times \Phi \times R_+)^n$ へ変換される。 $B_n$ を $n$ 番目の決定時刻までの可能な情報の集合と呼び、この可能な情報を通してのみ観測可能な履歴に依存する政策を情報政策(I-政策)と呼ぶ。I-政策は $\pi = \{\pi_0, \pi_1, \dots\}$ ,  $\pi_n \in Q(A|B_n)$ ,  $n \geq 0$ , である。すべての $b_n \in B_n$ に対して、 $\pi_n(f(\varphi_n)|b_n) = 1$ ,  $n \geq 0$ , なるボレル可測写像 $f: \Phi \rightarrow A$ が存在するとき、I-政策 $\pi$ は定常であるといい、 $f^\pi$ と表わす。任意のI-政策 $\pi$ に対して、 $\omega_n^\pi(\cdot | h_n) = \pi_n(\cdot | b_n^{h_n})$ ,  $n \geq 0$ , として政策 $\omega^\pi = \{\omega_0^\pi, \omega_1^\pi, \dots\}$ を定義する。ここで、 $b_n^{h_n}$ は $h_n \in H_n$ に対応する $B_n$ の要素である。このとき、 $\pi$ と $\omega^\pi$ は $A$ に同じ条件付確率を割り当てるので、すべてのI-政策は政策とみなすことができる。I-政策 $\pi$ と $\varphi, g, \hat{u}$ に対して $Q(S \times (A \times S \times R_+ \times M \times \Phi)^\infty | \Phi)$ の要素 $\bar{p}_\pi = \varphi \pi_0 g g \hat{u} \pi_1 g g \hat{u} \dots$ が定義できる。I-政策 $\pi$ を用いたときの期待合計割引利得を

$$(3.3) \quad J_\pi(\varphi) = \bar{E}_\pi \left[ \sum_{n=0}^{\infty} e^{-\alpha T_n} r(i_n, a_n) | \varphi \right], \quad \varphi \in \Phi,$$

と定義する。ここで、 $\bar{E}_\pi[\cdot | \varphi]$ は $\bar{p}_\pi\{\cdot | \varphi\}$ による条件付期待

値を表わす。政策  $\omega$  に対して、 $\mathcal{Q}(S \times (A \times S \times R_+ \times M \times \Phi)^\infty | \Phi)$  の要素  $\bar{\omega} = \omega_0 \bar{g} \hat{u} \omega_1 \bar{g} \hat{u} \dots$  を定義し、 $\bar{\omega} \{ \cdot | \varphi \}$  による条件付期待値を  $\bar{E}_\omega [ \cdot | \varphi ]$  と表わすと、 $J_\omega$  の定義において、 $E_\omega$  を  $\bar{E}_\omega$  で置きかえることができる。

$\bar{g} \in \mathcal{Q}(\Phi \times R_+ | \Phi \times A)$  を

$$(3.4) \quad \bar{g}(B | \varphi, a) = \sum_{i \in S} \sum_{j \in S} \sum_{m \in M} g(i, j) \times B_m(\varphi, a | i, a) g(m | j) \varphi(i),$$

$$B \in \mathcal{B}(\Phi \times R_+), \varphi \in \Phi, a \in A,$$

と定義する。ここで、

$$(3.5) \quad B_m(\varphi, a) = \{ t \mid (t, m) \in B(\varphi, a) \},$$

$$(3.6) \quad B(\varphi, a) = \{ (t, m) \mid (\hat{u}(\varphi, a, t, m), t) \in B \}.$$

また、 $\bar{r} \in M(\Phi \times A)$  を

$$(3.7) \quad \bar{r}(\varphi, a) = \sum_{i \in S} r(i, a) \varphi(i), \varphi \in \Phi, a \in A,$$

と定義する。POSMDP  $(S, M, A, g, \bar{g}, r, \alpha)$  は、以上のものから定義される SMDP  $(\Phi, A, \bar{g}, \bar{r}, \alpha)$  に同値に変換されることを示す。ここで、 $\Phi$  は状態の集合； $A$  は行動の集合； $\bar{g}$  は運動法則； $\bar{r}$  は利得関数； $\alpha$  は割引因子。SMDP に対する政策は I-政策と同じなので、 $\pi = \{ \pi_0, \pi_1, \dots \}$  と表わす。I-政策  $\pi$  と  $\bar{g}$  に対して、 $\mathcal{Q}((A \times \Phi \times R_+)^\infty | \Phi)$  の要素  $\bar{\pi} = \pi_0 \bar{g} \pi_1 \bar{g} \dots$

が定義できる。I-政策  $\pi$  を用いたときの期待合計割引利得を

$$(3.8) \quad I_{\pi}(\varphi) = E_{\pi} \left[ \sum_{n=0}^{\infty} e^{-\alpha T_n} r(\varphi_n, a_n) \mid \varphi \right], \quad \varphi \in \Phi,$$

と定義する。ここで、 $E_{\pi}[\cdot \mid \varphi]$  は  $\phi_{\pi}(\cdot \mid \varphi)$  による条件付期待値を表す。

Remark 2. Condition (1) が満たれるとすると、

$$(3.9) \quad \bar{g}(\Phi \times [0, \delta] \mid \varphi, a) \leq 1 - \varepsilon, \quad \varphi \in \Phi, a \in A,$$

が成り立つ。

Lemma 3.1. 任意の政策  $\omega$  と任意の I-政策  $\pi$  に対して、

$$(3.10) \quad J_{\omega}(\varphi) = \bar{E}_{\omega} \left[ \sum_{n=0}^{\infty} e^{-\alpha T_n} r(\varphi_n, a_n) \mid \varphi \right], \quad \varphi \in \Phi;$$

$$(3.11) \quad J_{\pi}(\varphi) = \bar{E}_{\pi} \left[ \sum_{n=0}^{\infty} e^{-\alpha T_n} r(\varphi_n, a_n) \mid \varphi \right], \quad \varphi \in \Phi,$$

が成り立つ。

Lemma 3.2. 任意の政策  $\omega$  と任意の I-政策  $\pi$  に対して、

$$(3.12) \quad \bar{g} \in \bar{\Phi}_{\omega}(\varphi_{n+1}, t_{n+1}) \mid (\varphi, a_0, \varphi_1, t_1, \dots, a_{n-1}, \varphi_n, t_n, a_n), \quad n \geq 0;$$

$$(3.13) \quad \bar{g} \in \bar{\Phi}_{\pi}(\varphi_{n+1}, t_{n+1}) \mid (\varphi, a_0, \varphi_1, t_1, \dots, a_{n-1}, \varphi_n, t_n, a_n), \quad n \geq 0,$$

が成り立つ。ここで、 $\bar{g} \in \bar{\Phi}_{\omega} x_2 \mid x_1$  ( $\bar{g} \in \bar{\Phi}_{\pi} x_2 \mid x_1$ ) は、 $\bar{g}$  は  $\bar{\Phi}_{\omega}$  ( $\bar{\Phi}_{\pi}$ ) の下で  $x_1$  が与えられたときの  $x_2$  の条件付分布であることを意味する。

Lemma 3.3. 各政策  $\omega$  に対して、 $J_{\pi\omega}(\varphi) = J_{\omega}(\varphi)$ ,  $\varphi \in \mathfrak{E}$ ,  
なる I-政策  $\pi^{\omega}$  が存在する。すなわち、I-政策で十分である。

以上の補題から次の定理を得る (cf. Wakuta [2,3])。

Theorem 3.4. 任意の I-政策  $\pi$  に対して、 $J_{\pi}(\varphi) = I_{\pi}(\varphi)$ ,  
 $\varphi \in \mathfrak{E}$ , が成り立つ。この意味で、 $\text{POSMDP}(S, M, A, g, \bar{g}, r, \alpha)$   
は  $\text{SMDP}(\mathfrak{E}, A, \bar{g}, \bar{r}, \alpha)$  と同値である。

#### §4. 最適定常 I-政策の存在

次の条件をおく。

Condition (2). (i) 各  $(i, j)$  に対して、 $k(j, \cdot | i, \cdot)$  は  
 $R_+ \times A$  上連続で、

$$(4.1) \quad k(j, t | i, a) \leq M_{ij}(t), \quad t \geq 0, \quad a \in A,$$

$$(4.2) \quad \int_0^{\infty} e^{-\alpha t} M_{ij}(t) dt < \infty,$$

なる可測関数  $M_{ij}(t)$  が存在する；

(ii)  $A$  はコンパクト距離空間である；

(iii)  $r(i, \cdot)$ ,  $i \in S$ , は  $A$  上の有界連続関数である。

いわゆる最適方程式は、

$$(4.3) \quad v(\varphi) = \max_{a \in A} \left\{ r(\varphi, a) + \int_{\mathfrak{E}} \int_0^{\infty} e^{-\alpha t} v(\varphi') d\bar{g}((\varphi', t) | \varphi, a) \right\},$$

$\varphi \in \mathfrak{E}$ ,

と表わされる。

$\bar{v}$  の定義から、この方程式は

$$(4.4) \quad v(\varphi) = \max_{a \in A} \left\{ T(\varphi, a) + \sum_{m \in M} \int_0^{\infty} e^{-\lambda t} v(\hat{u}(\varphi, a, t, m)) \right. \\ \left. \times \left[ \sum_{\hat{i} \in S} \sum_{\hat{j} \in S} k(\hat{j}, t | \hat{i}, a) g(m | \hat{j}) \varphi(\hat{i}) \right] dt \right\},$$

$\varphi \in \Phi,$

と表わすことが出来る。(4.3) あるいは (4.4) 式の右辺の  
カッコの項を  $v(\varphi, a)$  と表わし、

$$(4.5) \quad T_a v(\varphi) = v(\varphi, a), \quad \varphi \in \Phi, a \in A;$$

$$(4.6) \quad T v(\varphi) = \max_{a \in A} T_a v(\varphi), \quad \varphi \in \Phi;$$

$$(4.7) \quad T_f v(\varphi) = v(\varphi, f(\varphi)), \quad \varphi \in \Phi,$$

として、オペレーター  $T_a, T, T_f$  を導入する。ここで、 $f$  は  
 $\Phi$  から  $A$  へのボレル可測写像である。

Lemma 7.1. Conditions (1) と (2) が成り立つとする。

このとき、(i) (4.6) で定義されるオペレーター  $T$  は  $(\Phi)$  上  
の縮小写像である；(ii) (4.7) で定義されるオペレーター  $T_f$   
は  $M(\Phi)$  上の縮小写像である。

(証明) (i) を証明する。(ii) の証明は同様である。  $m \in M$   
を固定する。  $u_m(\varphi, a, t) = u(\varphi, a, t, m)$ ,  $v_m(\varphi, a, t) =$   
 $\sum_{\hat{i} \in S} \sum_{\hat{j} \in S} k(\hat{j}, t | \hat{i}, a) g(m | \hat{j}) \varphi(\hat{i})$ ,  $D_m = \{(\varphi, a, t) | v_m(\varphi, a, t) > 0\}$   
とおく。Condition (2) より  $u_m$  は  $D_m$  上連続で、 $v_m$  は  $\Phi \times A \times \mathbb{R}_+$   
上連続である。  $(\varphi, a, t)$  を  $D_m$  の境界点とし、  $(\varphi_n, a_n, t_n) \in D_m$ ,

$(\varphi_n, a_n, t_n) \rightarrow (\varphi, a, t) \ (n \rightarrow \infty)$  とする。このとき、 $v_m(\varphi_n, a_n, t_n) \rightarrow v_m(\varphi, a, t) = 0$ 。  $v$  は有界で、 $(\varphi, a, t) \notin D_m$  なら、 $v_m(\varphi, a, t) = 0$  なので、(4.4) の被積分関数は  $D_m$  の境界で連続で、したがって、全空間  $\Phi \times A \times R_+$  上で連続である。Condition (2) より、 $v(\varphi, a)$  は  $(\varphi, a)$  について連続であることが容易に示される。  $A$  はコンパクトなので、 $Tv(\varphi)$  は  $\varphi$  について連続である。したがって、 $T$  は  $C(\Phi)$  上のオペレーターである。(3.9) より、任意の  $v, v' \in C(\Phi)$  に対して、

$$(4.8) \quad \|Tv - Tv'\| \leq (1 - \varepsilon + \varepsilon e^{-\alpha\delta}) \|v - v'\|$$

が成り立ち、 $T$  は縮小写像である ( $\|\cdot\|$  はスーパールムを表わす)。以上で証明された。

Theorem 4.2. Conditions (1) と (2) が成り立つとする。このとき、 $I$ -政策  $\pi$  は、利得  $I(\pi)$  が最適方程式 (4.3) あるいは (4.4) を満たすときに限り最適である。

証明は Wakata [4] を参照のこと。

Theorem 4.3. Conditions (1) と (2) が成り立つとする。このとき、各  $\varphi$  に対して、 $T$  の不動点  $v^*$  をもつ (4.3) 式あるいは (4.4) 式の右辺を最大にする行動を選択する最適定常  $I$ -政策  $f$  が存在する。すなわち、 $f$  は次のようなものである：

$$\begin{aligned}
(4.9) \quad & \bar{V}(\varphi, f(\varphi)) + \sum_{m \in M} \int_0^{\infty} e^{-\alpha t} v^+(\hat{u}(\varphi, f(\varphi), t, m)) \\
& \times \left[ \sum_{i \in S} \sum_{\tilde{z} \in S} k(\tilde{z}, t | i, f(\varphi)) g(m | \tilde{z}) \varphi(i) \right] dt \\
= \max_{a \in A} & \left\{ \bar{V}(\varphi, a) + \sum_{m \in M} \int_0^{\infty} e^{-\alpha t} v^+(\hat{u}(\varphi, a, t, m)) \right. \\
& \left. \times \left[ \sum_{i \in S} \sum_{\tilde{z} \in S} k(\tilde{z}, t | i, a) g(m | \tilde{z}) \varphi(i) \right] dt \right\}
\end{aligned}$$

証明は、Theorem 4.2 と良く知った Dubins and Savage の selection theorem から直ちに成り立つ。

### References

- [1] Ross, S. (1970). Applied Probability Models with Optimization Applications. Holden-Day.
- [2] Wakata, K. (1981). Semi-Markov decision processes with incomplete state observation - Average cost criterion. J. Oper. Res. Soc. Japan, 24, 95-108.
- [3] \_\_\_\_\_ (1982). Semi-Markov decision processes with incomplete state observation - Discounted cost criterion. J. Oper. Res. Soc. Japan, 25, 351-361.
- [4] \_\_\_\_\_. Arbitrary state semi-Markov decision processes with unbounded rewards. To appear in Optimization.

- [5] White, C. (1976). Procedures for the solutions of a finite-horizon partially observed semi-Markov optimization problem. Oper. Res. 24, 348-358.