

# 不確実性の下でのマルコフ決定過程に対する区間ベイズ手法 (An Interval Bayesian Method for uncertain MDPs)

宮崎大学・教育文化学部 伊喜 哲一郎 (Tetsuichiro IKI)

Faculty of Education and Culture, Miyazaki University

神奈川大学・工学部 堀口 正之 (Masayuki Horiguchi)

Faculty of Engineering, Kanagawa University

千葉大学・理学部 安田 正實 (Masami YASUDA)

Faculty of Science, Chiba University

藏野 正美 (Masami KURANO)

## 1 はじめに

推移確率行列が未知のマルコフ決定過程 (Markov Decision Processes, MDPs) の解析は、最尤推定法を用いる場合 (cf. [2, 5, 7, 11]) とベイズ推定法を用いる場合 (cf. [5, 12, 15]) がある。ベイズ推定法においては、事前分布をいかに設定するかが一つの問題である。その設定において、柔軟性と融通性に富んだ頑健なモデルを構成することは現実問題への応用において重要である。

本論文では De Robertis and Hartigan[1] が提唱した事前測度区間による区間ベイズ法の考え方を適応して、推移確率行列が未知の MDPs の解析を試みる。そのために、未知の推移確率行列をある区間で推定した場合のモデルとして、区間推定 MDPs (Interval estimated MDPs) を定式化しその解析を行う。この解析結果を受けて、事前情報を区間ベイズ法 [1] にもとづく処理から得られた区間を用いたモデルとして、区間ベイズ MDPs (Interval Bayesian estimated MDPs) を構成する。

マルコフ連鎖の推移確率行列の区間ベイズ推定は、基本的には、多項分布の生起確率の区間推定に帰着されるので、これに関する計算法といいくつかの数値例を与える。Kurano et al[8, 9] で考察された “Controlled Markov set-chain model” は、推移確率行列を区間でとらえる考え方においては本論文と同じであるが、前者においては各期で推移確率行列が区間に変動することも可能な場合を取り扱っている。区間推定 MDPs では、全過程を通して推移確率行列は一定である場合を扱う。

## 2 記号と基本命題

ここでは、いくつかの記号と続く節で用いられる基本補題を与えておく。

$\mathbb{R}, \mathbb{R}^n, \mathbb{R}^{m \times n}$  をそれぞれ実数、 $n$  次元実列ベクトル、 $m \times n$  型実行列の全体を表す。 $\mathbb{R} = \mathbb{R}^{1 \times 1}, \mathbb{R}^n = \mathbb{R}^{n \times 1}$  とする。また、 $\mathbb{R}_+, \mathbb{R}_+^n, \mathbb{R}_+^{m \times n}$  はそれぞれ  $\mathbb{R}, \mathbb{R}^n, \mathbb{R}^{m \times n}$  の各成分が非負であるようなものの集合とする。 $\mathbb{R}^{m \times n}$  の半順序  $\preceq, \prec$  は次で定める:  $\mathbb{R}^{m \times n} \ni A = (a_{ij}), B = (b_{ij})$  に対して

$$(2.1) \quad A \preceq B (a_{ij} \leq b_{ij} (1 \leq i \leq m, 1 \leq j \leq n) のとき), A \prec B (A \preceq B かつ A \neq B のとき)$$

とする。 $\underline{A} \preceq \bar{A}$  なる  $\underline{A} = (\underline{a}_{ij}), \bar{A} = (\bar{a}_{ij}) \in \mathbb{R}_+^{m \times n}$  に対して区間  $\langle \underline{A}, \bar{A} \rangle$  を次で定める:

$$(2.2) \quad \langle \underline{A}, \bar{A} \rangle = \{Q = (q_{ij}) \in \mathbb{R}_+^{m \times n} \mid \underline{a}_{ij} \leq q_{ij} \leq \bar{a}_{ij}, q_{ij} \geq 0, \sum_{j=1}^n q_{ij} = 1 (1 \leq i \leq m, 1 \leq j \leq n)\}.$$

$n \times n$  型の確率行列の区間集合全体を  $\mathcal{M}_n = \{\langle Q, \bar{Q} \rangle \mid \langle Q, \bar{Q} \rangle \neq \emptyset, Q \preceq \bar{Q}, Q, \bar{Q} \in \mathbb{R}_+^{n \times n}\}$  で表す。 $\mathcal{M}_n \ni Q_1, Q_2$  に対する積  $Q_1 Q_2$  を  $Q_1 Q_2 = \{Q_1 Q_2 \mid Q_1 \in \mathcal{M}_n, Q_2 \in \mathcal{M}_n\}$  と定める。また、 $Q \in \mathcal{M}_n$  に対する多重積は逐次的に定義される:  $Q^k = Q^{k-1} Q$  ( $k \geq 2$ )。

$C(\mathbb{R}_+)$  を  $\mathbb{R}_+$  の有界閉区間の全体とする。また、 $C(\mathbb{R}_+)^n$  を  $C(\mathbb{R}_+)$  の要素を成分に持つ  $n$  次元列ベクトルの全体とする:  $C(\mathbb{R}_+)^n = \{D = (D_1, D_2, \dots, D_n)' \mid D_i \in C(\mathbb{R}_+) (1 \leq i \leq n)\}$ 。ただし、 $\mathbf{d}'$  はベクトル  $\mathbf{d}$  の転置を表す。 $C(\mathbb{R}_+)^n$  上の算法 (加法、スカラ倍) は次で定める:

$D = (D_1, D_2, \dots, D_n)', E = (E_1, E_2, \dots, E_n)' \in C(\mathbb{R}_+)^n, h \in \mathbb{R}_+^n, \lambda \in \mathbb{R}_+$  に対して,

$$(2.3) \quad D + E = \{d + e \mid d \in D, e \in E\}, h + D = \{h + d \mid d \in D\}, \lambda D = \{\lambda d \mid d \in D\}.$$

$D = ([d_1, \bar{d}_1], [d_2, \bar{d}_2], \dots, [d_n, \bar{d}_n])' \in C(\mathbb{R}_+)^n$  を  $D = [\underline{d}, \bar{d}]$  と記す。ただし、 $\underline{d} = (d_1, d_2, \dots, d_n) \in \mathbb{R}_+^n, \bar{d} = (\bar{d}_1, \bar{d}_2, \dots, \bar{d}_n) \in \mathbb{R}_+^n$  とする。 $D = (D_1, D_2, \dots, D_n)' \in C(\mathbb{R}_+)^n$  と部分集合  $G \subset \mathbb{R}_+^{1 \times n}$  に対し

て, その積  $GD$  を  $GD = \{gd | g = (g_1, g_2, \dots, g_n) \in G, d = (d_1, d_2, \dots, d_n)' \in D, d_i \in D_i (1 \leq i \leq n)\}$  と定める.

次が成り立つ.

**Lemma 2.1.** ([4, 8]) (i) 任意の  $\mathcal{Q} \in \mathcal{M}_n$  は  $n \times n$  次元ベクトル空間  $\mathbb{R}^{n \times n}$  の凸多面体である. (ii) コンパクト凸部分集合  $G \subset \mathbb{R}_+^{1 \times n}$  と  $D = (D_1, D_2, \dots, D_n) \in C(\mathbb{R}_+)^n$  に対して  $GD \in C(\mathbb{R}_+)$  である.

$C(\mathbb{R}_+)$  上の半順序  $\preceq, \prec$  を次で定める:  $[c_1, c_2], [d_1, d_2] \in C(\mathbb{R}_+)$  に対して

$$(2.4) \quad \begin{cases} [c_1, c_2] \preceq [d_1, d_2] & (c_i \leq d_i (i = 1, 2) \text{ のとき}) \\ [c_1, c_2] \prec [d_1, d_2] & ([c_1, c_2] \preceq [d_1, d_2] \text{ かつ } [c_1, c_2] \neq [d_1, d_2] \text{ のとき}) \end{cases}$$

とする.  $C(\mathbb{R}_+)^n$  上の半順序  $\preceq, \prec$  は  $C(\mathbb{R}_+)$  上の半順序を用いて次により定める:  $\mathbf{v} = (v_1, v_2, \dots, v_n)', \mathbf{w} = (w_1, w_2, \dots, w_n)' \in C(\mathbb{R}_+)^n$  に対して

$$(2.5) \quad \mathbf{v} \preceq \mathbf{w} (v_i \leq w_i (1 \leq i \leq n) \text{ のとき}), \mathbf{v} \prec \mathbf{w} (\mathbf{v} \prec \mathbf{w} \text{ かつ } \mathbf{v} \neq \mathbf{w} \text{ のとき})$$

$\mathbb{R}_+^n$  の 2 つの有界閉集合  $D_1, D_2$  の距離としてハウスドルフ距離  $\rho$  を考える:

$$(2.6) \quad \rho(D_1, D_2) = \max\{\sup_{x \in D_1} \inf_{y \in D_2} \|x - y\|, \sup_{y \in D_2} \inf_{x \in D_1} \|x - y\|\}.$$

ただし,  $\|\cdot\|$  は  $\mathbb{R}^n$  におけるユークリッド距離とする.

次に, 次節以降の議論の準備として有限状態マルコフ決定過程について述べる. ある決定過程の状態空間を  $S = \{1, 2, \dots, n\}$ , 行動空間を  $A = \{1, 2, \dots, k\}$  とする. 次の集合を定義する:

$$(2.7) \quad P(S) := \{p = (p_1, p_2, \dots, p_n) \in \mathbb{R}_+^n \mid \sum_{i \in S} p_i = 1\},$$

$$(2.8) \quad P(S|S) := \{q = (q_{ij} : i, j \in S) \in \mathbb{R}_+^{n \times n} \mid \sum_{j \in S} q_{ij} = 1 (i \in S)\},$$

$$(2.9) \quad P(S|S \times A) := \{Q = (q_{ij}(a) : i, j \in S, a \in A) \in \mathbb{R}_+^{kn \times n} \mid q_{i.}(a) \in P(s) (i \in S, a \in A)\}.$$

有限集合  $D$  上の非負実数値関数の全体を  $B_+(D)$  で表す.  $D$  が有限集合のとき  $B_+(D)$  と  $\mathbb{R}_+^n$  を同一視する. ただし  $n = |D|$  であるとする.  $Q = (q_{ij}(a)) \in P(S|S \times A)$  と  $\mathbf{r} = (r(i, a)) \in B_+(S \times A)$  に対して, 通常のマルコフ決定過程 MDPs  $\{S, A, Q, \mathbf{r}\}$  を考え (cf. [13]), ここでは簡単のために確定的 (deterministic) で定常 (stationary) な政策のみを考える.  $S$  から  $A$  への写像  $f$  の全体を  $F$  で表す. 任意の  $f \in F$  に対して, 割引率  $\beta (0 < \beta < 1)$  によって割り引かれた総期待利得ベクトル  $\phi(f|Q) \in \mathbb{R}_+^n$  を確率行列  $Q \in P(S|S \times A)$  の関数として次で定める:

$$(2.10) \quad \phi(f|Q) = \sum_{t=0}^{\infty} (\beta Q(f))^t \mathbf{r}(f),$$

ただし,  $\mathbf{r}(f) = (r(1, f(1)), r(2, f(2)), \dots, r(n, f(n)))' \in \mathbb{R}_+^n$ ,  $Q(f) = (q_{ij}(f(i))) \in P(S|S)$ . 各  $f \in F$  に対して写像  $L(f) : \mathbb{R}_+^n \rightarrow \mathbb{R}_+^n$  を次で定める:

$$(2.11) \quad L(f)\mathbf{x} = \mathbf{r}(f) + \beta Q(f)\mathbf{x}, \quad \mathbf{x} = (x_1, x_2, \dots, x_n)' \in \mathbb{R}_+^n.$$

このとき, 次の基本補題が知られている.

**Lemma 2.2.** (cf. [13]) (i)  $L(f)$  は単調増加および縮小写像である. すなわち,  $\mathbf{x} \leq \mathbf{x}'$  ならば  $L(f)\mathbf{x} \leq L(f)\mathbf{x}'$  (componentwise),  $\|L(f)\mathbf{x} - L(f)\mathbf{x}'\| \leq \beta \|\mathbf{x} - \mathbf{x}'\|$  ( $\mathbf{x}, \mathbf{x}' \in \mathbb{R}_+^n$ ) が成り立つ. ただし,  $\|\cdot\|$  は sup ノルムとする. (ii)  $\phi(f|Q)$  は  $L(f)$  の唯一の不動点である. すなわち任意の  $\mathbf{x} \in \mathbb{R}_+^n$  に対して  $L(f)^t \mathbf{x} \rightarrow \phi(f|Q)$  ( $t \rightarrow \infty$ ) が成り立つ.

### 3 区間推定 MDPs とパレート最適

本節では、MDPs $\{S, A, Q, \mathbf{r}\}$  の推移確率行列  $Q$  を区間  $\mathcal{Q} = \langle \underline{Q}, \bar{Q} \rangle$  で推定した場合を考察する。ただし、

$$(3.1) \quad \underline{Q} = (\underline{q}_{ij}(a) : i, j \in S, a \in A) \in \mathbb{R}_+^{kn \times n}, \bar{Q} = (\bar{q}_{ij}(a) : i, j \in S, a \in A) \in \mathbb{R}_+^{kn \times n},$$

$$(3.2) \quad \mathcal{Q} = \langle \underline{Q}, \bar{Q} \rangle = \{Q \in P(S|S \times A) \mid \underline{Q} \leqq Q \leqq \bar{Q}\}$$

とする。推移確率行列  $Q$  を  $\mathcal{Q} = \langle \underline{Q}, \bar{Q} \rangle$  で推定した決定モデルを区間推定 MDPs $\{\mathcal{Q}\}$ (Interval estimated MDPs $\{\mathcal{Q}\}$ ) と呼ぶ。以下、区間推定 MDPs の利得関数を定義しその最適化について議論する。

$f \in F$  に対する割引された総期待-集合ベクトル  $\phi(f|\mathcal{Q})$  を次で定める:

$$(3.3) \quad \phi(f|\mathcal{Q}) = \{\phi(f|Q) \mid Q \in \mathcal{Q}\} \subset \mathbb{R}_+^n.$$

ただし、 $\phi(f|Q)$  は式 (2.10) で与えられている。ここで、 $\phi(f|\mathcal{Q}) \in C(\mathbb{R}_+)^n$  であることを示そう。 $\mathcal{L}$  を  $C(\mathbb{R}_+)^n$  から  $C(\mathbb{R}_+)^n$  への写像で次のように定める:

$$(3.4) \quad \mathcal{L}(f)\mathbf{v} = \mathbf{r}(f) + \beta \mathcal{Q}(f)\mathbf{v}, \quad \mathbf{v} \in C(\mathbb{R}_+)^n,$$

ただし、式 (3.4)において  $\mathcal{Q}(f) = \langle \underline{Q}(f), \bar{Q}(f) \rangle$ ,  $\underline{Q}(f) = (\underline{q}_{ij}(f(i))) \in \mathbb{R}_+^{n \times n}$ ,  $\bar{Q}(f) = (\bar{q}_{ij}(f(i))) \in \mathbb{R}_+^{n \times n}$  である。Lemma 2.1 により  $\mathcal{L}(f)\mathbf{v} \in C(\mathbb{R}_+)^n$  ( $\mathbf{v} \in C(\mathbb{R}_+)^n$ ) であることが示されていることに注意する。さらに、 $\underline{L}(f) : \mathbb{R}_+^n \rightarrow \mathbb{R}_+^n$ ,  $\bar{L}(f) : \mathbb{R}_+^n \rightarrow \mathbb{R}_+^n$  を次で定める:  $\mathbf{x} = (x_1, x_2, \dots, x_n)' \in \mathbb{R}_+^n$  に対して

$$(3.5) \quad \underline{L}(f)\mathbf{x} = \mathbf{r}(f) + \beta \min_{Q \in \mathcal{Q}(f)} Q\mathbf{x}, \quad \bar{L}(f)\mathbf{x} = \mathbf{r}(f) + \beta \max_{Q \in \mathcal{Q}(f)} Q\mathbf{x}.$$

このとき、次が成り立つ。

**Lemma 3.1.** 任意の  $f \in F$  に対して、次が成り立つ: (i)  $\mathcal{L}(f)$  は単調増加かつ縮小写像である。 (ii)  $\underline{L}(f), \bar{L}(f)$  は、ともに単調増加かつ sup ノルムに関して縮小写像である。

*Proof.* [8] の定理 3.1 を参照。 ■

Lemma 2.2 と Lemma 3.1 を適用して次を得る。

**Theorem 3.1.** 任意の  $f \in F$  に対して次が成り立つ: (i)  $\phi(f|\mathcal{Q}) \in C(\mathbb{R}_+)^n$  かつ  $\phi(f|\mathcal{Q})$  は  $\mathcal{L}(f)$  の唯一の不動点である。さらに、任意の  $\mathbf{v} \in C(\mathbb{R}_+)^n$  に対して  $\mathcal{L}(f)^\ell \mathbf{v} \rightarrow \phi(f|\mathcal{Q})$  ( $\ell \rightarrow \infty$ )。 (ii)  $\phi(f|\mathcal{Q}) = [\underline{\phi}(f), \bar{\phi}(f)]$  とすると、 $\underline{\phi}(f), \bar{\phi}(f)$  はそれぞれ  $\underline{L}(f), \bar{L}(f)$  の唯一の不動点である。

*Proof.* 任意の  $Q \in \mathcal{Q}$  に対して、 $\phi(f|Q) = \mathbf{r}(f) + \beta Q(f)\phi(f|Q) \leqq \bar{L}(f)\phi(f|Q)$  これより、 $\phi(f|Q) \leqq \bar{L}(f)^\ell \phi(f|Q) \rightarrow \bar{\phi}(f)$  ( $\ell \rightarrow \infty$ )。同様にして、 $\phi(f|Q) \geqq \underline{L}(f)^\ell \phi(f|Q) \rightarrow \underline{\phi}(f)$  ( $\ell \rightarrow \infty$ )。故に、 $\underline{\phi}(f) \leqq \phi(f|Q) \leqq \bar{\phi}(f)$  を得る。明らかに、 $\underline{\phi}(f), \bar{\phi}(f) \in \phi(f|\mathcal{Q})$  かつ  $\phi(f|Q)$  は  $Q \in \mathcal{Q}$  に関して連続(cf. [14])であるから  $\phi(f|\mathcal{Q}) = [\underline{\phi}(f), \bar{\phi}(f)]$  が成り立つ。これで (ii) が示された。

$\mathcal{L}(f)$  の不動点を  $\mathbf{u}(f) \in C(\mathbb{R}_+)^n$  とする。任意の  $\mathbf{v} = [\underline{v}, \bar{v}] \in C(\mathbb{R}_+)^n$  に対して次が成り立つ:  $\mathcal{L}(f)\mathbf{v} = [\underline{L}v, \bar{L}v]$ 。故に、 $\ell \geqq 1$  に対して  $\mathcal{L}(f)^{\ell \ell} \mathbf{v} = [\underline{L}^\ell v, \bar{L}^\ell \bar{v}]$ 。これより、 $\ell \rightarrow \infty$  とすることで  $\mathbf{u}(f) = [\underline{\phi}(f), \bar{\phi}(f)]$  を得る。(ii) より  $\mathbf{u}(f) = \phi(f|\mathcal{Q})$  となり (i) が示された。 ■

$f^* \in F$  がパレート最適であるとは、 $\phi(f^*|\mathcal{Q}) \prec \phi(f|\mathcal{Q})$  なる  $f \in F$  が存在しない場合を言う。

**Lemma 3.2.**  $f, g \in F$  に対して、 $\phi(f|\mathcal{Q}) \prec \mathcal{L}(g)\phi(f|\mathcal{Q})$  ならば  $\phi(f|\mathcal{Q}) \prec \phi(g|\mathcal{Q})$ 。

*Proof.*  $\mathcal{L}(g)$  の単調性と Theorem 3.1 から  $\phi(f|\mathcal{Q}) \prec \mathcal{L}(g)\phi(f|\mathcal{Q}) \prec \mathcal{L}(g)\mathcal{L}(g)\phi(f|\mathcal{Q}) \prec \cdots \prec (\mathcal{L}(g))^n \phi(f|\mathcal{Q}) \rightarrow \phi(g|\mathcal{Q})$  ( $n \rightarrow \infty$ )。従って、 $\phi(f|\mathcal{Q}) \prec \phi(g|\mathcal{Q})$  が示された。 ■

$D \subset C(\mathbb{R}_+)^n$  に対して点  $\mathbf{v} \in D$  が  $D$  の有効点 (efficient point) であるとは、 $\mathbf{v} \prec \mathbf{u}$  なる  $\mathbf{u} \in D$  が存在していない場合を言う。 $D$  の有効点の全体を  $\text{eff}(D)$  で表す。式 (3.1) の  $\underline{Q}, \bar{Q}$  の成分ベクトル  $\underline{Q}_{i,a} = (\underline{q}_{i1}(a), \underline{q}_{i2}(a), \dots, \underline{q}_{in}(a)), \bar{Q}_{i,a} = (\bar{q}_{i1}(a), \bar{q}_{i2}(a), \dots, \bar{q}_{in}(a))$  に対して  $\mathcal{Q}_{i,a} = \langle \underline{Q}_{i,a}, \bar{Q}_{i,a} \rangle$  ( $i \in S, a \in A$ ) とする。 $\mathbf{u} \in C(\mathbb{R}_+)^n$  に対して次を定める:

$$(3.6) \quad \mathcal{L}(\mathbf{u}) := (\mathcal{L}(\mathbf{u})_1, \mathcal{L}(\mathbf{u})_2, \dots, \mathcal{L}(\mathbf{u})_n)',$$

ただし、 $\mathcal{L}(\mathbf{u})_i := \text{eff}(\{r(i, a) + \beta \mathcal{Q}_{i,a} \mathbf{u} \mid a \in A\})$  ( $i \in S$ ) である。

このとき、Lemma 3.2 を用いて次が示される。

**Theorem 3.2.**  $f^*$  がパレート最適であるための必要十分条件は,  $\phi(f^*|\mathcal{Q})$  が次の最適包含式の最大解となることである.

$$(3.7) \quad \mathbf{u} \in \mathcal{L}(\mathbf{u}), \mathbf{u} \in C(\mathbb{R}_+)^n.$$

*Proof.* ( $\Rightarrow$ )  $f^* \in F$  を Pareto-optimal とする. このとき, Theorem 3.1 から  $\phi(f^*|\mathcal{Q})$  は  $\mathcal{L}(f^*)$  の不動点である. よって,  $\phi(f^*|\mathcal{Q}) \in \mathcal{L}(\phi(f^*|\mathcal{Q}))$  である. ここで, もし, ある  $\mathbf{u} \in C(\mathbb{R}_+)^n$  が存在して  $\mathbf{u} \in \mathcal{L}(\mathbf{u})$ かつ  $\phi(f^*|\mathcal{Q}) \prec \mathbf{u}$  であるものが存在したとする. すなわち,  $\exists g \in F, \exists i_0 \in S$  s.t.  $\phi(f^*|\mathcal{Q})_{i_0} \prec \mathbf{u}_{i_0} = r(i_0, g(i_0)) + \beta \mathcal{Q}_{i_0, g(i_0)} \mathbf{u} = \phi(g|\mathcal{Q})_{i_0}, \phi(f^*|\mathcal{Q})_i \preceq \mathbf{u}_i = r(i, g(i)) + \beta \mathcal{Q}_{i, g(i)} \mathbf{u} = \phi(g|\mathcal{Q})_i$  ( $i \neq i_0, i \in S$ ). ただし,  $\mathbf{a}_i$  は  $\mathbf{a} \in C(\mathbb{R}_+)^n$  の第  $i$  成分を表す. これは,  $f^*$  が Pareto-optimal であることに矛盾する.  
( $\Leftarrow$ )  $\phi(f^*|\mathcal{Q})$  を  $\mathbf{u} \in \mathcal{L}(\mathbf{u})$  の最大解であるが Pareto-optimal でないとする. このとき,  $\exists g \in F$  s.t.  $\phi(f^*|\mathcal{Q}) \prec \phi(g|\mathcal{Q})$  である. 特に,  $\exists i \in S$  s.t.  $\phi(f^*|\mathcal{Q})_i \prec \phi(g|\mathcal{Q})_i$  である. 一般に,  $\phi(f|\mathcal{Q})_i \in C(\mathbb{R}_+)$  であって,  $f \in F$  は高々有限個であるから  $\phi(g|\mathcal{Q})_i \preceq \phi(\bar{g}|\mathcal{Q})_i$  ( $i \in S$ ) となる  $\bar{g}$  が存在する. すなわち  $\phi(f^*|\mathcal{Q}) \prec \phi(g|\mathcal{Q}) \preceq \phi(\bar{g}|\mathcal{Q})$  が成り立つ. ここで,  $\phi(\bar{g}|\mathcal{Q}) \notin \mathcal{L}(\phi(\bar{g}|\mathcal{Q}))$  であると仮定すると, 有効点の定義から  $\exists i_0 \in S$  と  $\exists a_0 \in A$  に対して

$$(3.8) \quad \phi(\bar{g}|\mathcal{Q})_{i_0} \prec r(i_0, a_0) + \beta \mathcal{Q}_{i_0, a_0} \phi(\bar{g}|\mathcal{Q})_{i_0}$$

が成り立つ. ここで  $f^{(1)}(i) = a_0$  (if  $i = i_0$ ),  $\bar{g}(i)$  (if  $i \neq i_0, i \in S$ ) とすれば, Lemma 3.2 から  $\phi(\bar{g}|\mathcal{Q}) \prec \phi(f^{(1)}|\mathcal{Q})$  を得る.  $\phi(f^*|\mathcal{Q}) \prec \phi(f^{(1)}|\mathcal{Q})$  であって,  $\phi(f^{(1)}|\mathcal{Q}) \notin \mathcal{L}(\phi(f^{(1)}|\mathcal{Q}))$  であれば式 (3.8) と同様にして  $\exists f^{(2)} \in F$  s.t.  $\phi(f^{(1)}|\mathcal{Q}) \prec \phi(f^{(2)}|\mathcal{Q})$ .  $f \in F$  は高々有限個だから,  $\exists f^{(k)} \in F$  s.t.  $\phi(f^*|\mathcal{Q}) \prec \phi(f^{(1)}|\mathcal{Q}) \prec \cdots \prec \phi(f^{(k)}|\mathcal{Q})$  and  $\phi(f^{(k)}|\mathcal{Q}) \in \mathcal{L}(\phi(f^{(k)}|\mathcal{Q}))$  が成り立つ. これは,  $\phi(f^*|\mathcal{Q})$  が  $\mathbf{u} \in \mathcal{L}(\mathbf{u})$  の最大解であることに矛盾する. ■

## 4 ディリクレ分布

マルコフ連鎖の推移確率行列の区間ベイズ推定は, 行列の行成分に着目すれば, 多項分布の生起確率の区間推定に帰着される. そこで, 次節以降に用いられる区間ベイズ法による推移確率の事前・事後解析のためにディリクレ分布(多次元ベータ分布)に関するいくつかの性質を示す.

ガンマ関数  $\Gamma(x)$  ( $x > 0$ ) とベータ関数  $B(x, y)$  ( $x, y > 0$ ) をそれぞれ次のように表すことにする.

$$\Gamma(x) = \int_0^\infty t^{x-1} e^{-t} dt \quad (x > 0), B(x, y) = \int_0^1 t^{x-1} (1-t)^{y-1} dt \quad (x, y > 0).$$

$k$ -変数ディリクレ分布の p.d.f. を次のように定義する:

$$(4.1) \quad f(x_1, \dots, x_k) = \frac{\Gamma(\nu_1 + \cdots + \nu_{k+1})}{\Gamma(\nu_1) \cdots \Gamma(\nu_{k+1})} x_1^{\nu_1-1} \cdots x_k^{\nu_k-1} (1 - x_1 - x_2 - \cdots - x_k)^{\nu_{k+1}-1}.$$

但し,  $x_1, \dots, x_k$  は  $k$  次元多面体  $S_k := \{(x_1, \dots, x_k) : x_i \geq 0, i = 1, \dots, k, \sum_{i=1}^k x_i \leq 1\}$  の各成分であり,  $f$  は  $S_k$  上の点以外では 0,  $\nu_i \in \mathbb{R}$  は  $\nu_i > 0$  ( $i = 1, 2, \dots, k+1$ ) とする.

$$(4.2) \quad \tilde{D}(\nu_1, \dots, \nu_k; \nu_{k+1}) = \int \cdots \int_{S_k} f(x_1, \dots, x_k) dx_1 \cdots dx_k$$

と表す. ディリクレ積分, すなわち, ディリクレ分布の定数係数を除いた被積分関数部分に関して

$$(4.3) \quad \begin{aligned} & D(\nu_1, \nu_2, \dots, \nu_k; \nu_{k+1}) \\ &:= \int \cdots \int_{S_k} x_1^{\nu_1-1} x_2^{\nu_2-1} \cdots x_k^{\nu_k-1} (1 - x_1 - x_2 - \cdots - x_k)^{\nu_{k+1}-1} dx_1 dx_2 \cdots dx_k \\ &= \frac{\Gamma(\nu_1) \cdots \Gamma(\nu_{k+1})}{\Gamma(\nu_1 + \cdots + \nu_{k+1})} = \prod_{n=1}^{k+1} B\left(\nu_n, \sum_{l=n+1}^{k+1} \nu_l\right) \end{aligned}$$

を得る.

$0 < \lambda < 1$  に対して,

$$(4.4) \quad D(\nu_1, \dots, \nu_k; \nu_{k+1} | \lambda) := \int \cdots \int_{S_k \cap \{0 < x_1 \leq \lambda\}} x_1^{\nu_1-1} \cdots x_k^{\nu_k-1} (1 - x_1 - \cdots - x_n)^{\nu_{k+1}-1} dx_1 \cdots dx_k \quad (k \geq 1)$$

とする. 特に  $B(\alpha, \beta | \lambda) := D(\alpha; \beta | \lambda)$  ( $\alpha, \beta > 0$ ) と表すとき,

$$(4.5) \quad D(\nu_1, \dots, \nu_k; \nu_{k+1} | \lambda) = B(\nu_1, \nu_2 + \cdots + \nu_{k+1} | \lambda) B(\nu_2, \nu_3 + \cdots + \nu_{k+1} | \lambda) \cdots B(\nu_k, \nu_{k+1} | \lambda)$$

が成り立つ. ここで,  $m, n$  を正の整数とするとき

$$(4.6) \quad B(m, n | \lambda) = \int_0^\lambda x^{m-1} (1-x)^{n-1} dx = \sum_{i=0}^{n-1} \binom{n-1}{i} (-1)^i \lambda^{m+i} \frac{1}{m+i} \quad (m, n > 0)$$

を得る.

## 5 区間ベイズ法による事前・事後解析

ここでは, De Robertis & Hartigan[1] による事前測度区間を用いた区間ベイズ法を定常マルコフ決定過程の推移確率行列の区間推定へ適用し, 区間推定 MDPs を構成する事後区間について考察する.

$P(S) = P_n = \{p = (p_1, p_2, \dots, p_n) | p_i \geq 0, \sum_{i=1}^n p_i = 1\}$  とおく.  $L(\cdot)$  を  $P_n$  上のルベーグ測度 (lower bound measure),  $U(\cdot) := kL(\cdot)$  (upper bound measure) を測度  $L$  の  $k(k > 0)$  に関する比例測度 (proportional measure) とし, 事前測度区間を  $[L, kL] = [dp, kdp]$  とする. データ  $\sigma = (\sigma_1, \sigma_2, \dots, \sigma_n)$  はある状態における  $\hat{\sigma} := \sum_{k=1}^n \sigma_k$  回の独立試行実験でそれぞれ状態  $i$  が  $\sigma_i$  回起きたことを表す. 状態  $i$  の生起確率が  $p_i$  であるとき,  $p = (p_1, \dots, p_n) \in P_n$  に対するデータ  $\sigma$  の p.d.f. は多項分布で表されて

$$(5.1) \quad f(\sigma_1, \sigma_2, \dots, \sigma_n | p) = \frac{(\sigma_1 + \cdots + \sigma_n)!}{\sigma_1! \cdots \sigma_n!} p_1^{\sigma_1} p_2^{\sigma_2} \cdots p_n^{\sigma_n}$$

となる.

データ  $\sigma$  における事後測度区間を  $[L_\sigma, U_\sigma] = [L_\sigma, kL_\sigma]$  とする. 次の期に状態  $i$  へ推移する確率  $p_i$  のうち, まず,  $p_1$  に関する事後測度区間

$$\left\{ \frac{\int_{P_n} p_1 Q(dp)}{\int_{P_n} Q(dp)} \middle| L_\sigma \leq Q \leq U_\sigma \right\}$$

について調べる. 論文 [1] から, 上の事後測度区間  $[\underline{\lambda}, \bar{\lambda}]$  は次の方程式の一意の解である.

$$(5.2) \quad U_\sigma(p_1 - \underline{\lambda})^- + L_\sigma(p_1 - \underline{\lambda})^+ = 0$$

$$(5.3) \quad U_\sigma(p_1 - \bar{\lambda})^+ + L_\sigma(p_1 - \bar{\lambda})^- = 0$$

ただし,  $x^+ = \max\{0, x\}$ ,  $x^- = x - x^+ = \min\{0, x\}$  である.

$\hat{\sigma} = \sigma_1 + \sigma_2 + \cdots + \sigma_n$ ,  $s = \sigma_1 + 1$ ,  $t = \hat{\sigma} - \sigma_1 + n - 1$  とおくと, 式 (5.2) と式 (5.3) は結局,

$$(5.4) \quad K(s, t, \lambda) := \left( \frac{s}{s+t} - \lambda \right) B(s, t) + (k-1)(B(s+1, t | \lambda) - \lambda B(s, t | \lambda)) = 0$$

$$(5.5) \quad G(s, t, \lambda) := k \left( \frac{s}{s+t} - \lambda \right) B(s, t) - (k-1)(B(s+1, t | \lambda) - \lambda B(s, t | \lambda)) = 0$$

の方程式の解として表される. いずれも  $(\hat{\sigma} + n)$  次多項式による方程式の解である.

**Theorem 5.1.** データ  $\sigma = (\sigma_1, \sigma_2, \dots, \sigma_n)$ ,  $\hat{\sigma} = \sum_{i=1}^n \sigma_i$  とする. 事前測度区間を  $[L, kL]$  とするとき,  $p = (p_1, p_2, \dots, p_n)$  の  $p_i$  についての事後測度区間  $[\underline{\lambda}, \bar{\lambda}]$  は次のそれぞれの方程式の一意の実数解である.

$$K(\sigma_i + 1, \hat{\sigma} + n - \sigma_i - 1, \lambda) = 0, G(\sigma_i + 1, \hat{\sigma} + n - \sigma_i - 1, \lambda) = 0.$$

## 6 A numerical experiment

前節までの多項分布に関して状態の個数  $n = 3$  のときを考える。 $P_3 = \{p = (p_1, p_2, p_3) | \sum_{i=1}^3 p_i = 1, p_i \geq 1, i = 1, 2, 3\}$  とおき,  $k = 2$  とする, すなわち事前測度区間を  $[L, 2L]$  とする。ある決まった状態から  $\hat{\sigma} = 6$  回の試行がなされ, 6回中, 状態 1 に 3 回, 状態 2 に 1 回, 状態 3 に 2 回推移したとする。よって,  $\sigma_1 = 3, \sigma_2 = 1, \sigma_3 = 2$  であり,  $\hat{\sigma} = \sigma_1 + \sigma_2 + \sigma_3 = 6, s = \sigma_1 + 1 = 4, t = \sigma_2 + \sigma_3 + (n - 1) = 5$  のデータが得られているとする。Theorem 5.1 より,  $\lambda$  に関する 9 次方程式  $8 - 18\lambda + \lambda^5(126 - 336\lambda + 360\lambda^2 - 180\lambda^3 + 35\lambda^4) = 0$  を解いて, 解  $\lambda \approx 0.489$ を得る。また,  $\lambda$  に関する方程式  $4 - 9\lambda - \lambda^5(126 - 336\lambda + 360\lambda^2 - 180\lambda^3 + 35\lambda^4) = 0$  を解いて, 解として  $\lambda \approx 0.400$ を得る。よって  $p_1$  の事後測度区間は  $[0.400, 0.489]$  と考えられる。

$k = 1$ , すなわち, 事前測度区間としてルベーグ測度を考えたとき, 事後測度区間を求める方程式から  $p_i = [\underline{p}_i, \bar{p}_i] = \frac{\sigma_i+1}{\hat{\sigma}+n}$  と 1 点で表される。これは, 一様事前分布を考えたときの観測値  $(\sigma_1, \sigma_2, \sigma_3)$  によるディリクレ分布(多次元ベータ分布)のパラメータ  $p_i$  の事後分布での期待値に等しい。

さらに, 数値実験を行い事後測度区間をもとにした Markov set-chain の問題を解いてみる(cf. [6])。状態数  $n = 3, S = \{1, 2, 3\}$ , policy は固定(deterministic stationary policy)として初期状態  $x_1 = 1$  から第 20 期の状態  $x_{20}$  を観測するまでのうち, それぞれの状態から次の期に推移した頻度を調べたところ, Table 6.1 左上のような行列であったとする。例えば, 状態 2 からの推移では, この行列の第 2 行目を見て, 6 回の試行実験で次の期にそれぞれ状態 1 に  $\sigma_1 = 1$  回, 状態 2 に  $\sigma_2 = 3$  回, 状態 3 に  $\sigma_3 = 3$  回の推移を観測したとする。

各状態  $i$  における  $p_{i1}, p_{i2}, p_{i3}$  の事後測度区間は, 本文の Theorem 5.1 から以下のように得られる(Table 6.1).  $\sigma_1, \sigma_2, \sigma_3$  はそれぞれ状態  $i$  での観測値(推移回数)とする。

Table 6.1: Intervals of posterior measures

| 状態 1, $\hat{\sigma} = 6$ (実験回数), $\sigma_1 = 3, \sigma_2 = 1, \sigma_3 = 2$ のとき: |   |   |
|--|---|---|
| 状態の観測度数: $\begin{pmatrix} 3 & 1 & 2 \\ 1 & 3 & 2 \\ 1 & 2 & 4 \end{pmatrix}$     | $\hat{p}_{11} = [\underline{p}_{11}, \bar{p}_{11}]$ | $\hat{p}_{12} = [\underline{p}_{12}, \bar{p}_{12}]$ |
|  | $[0.400, 0.489]$                                    | $[0.187, 0.260]$                                    |
| 状態 2, $\hat{\sigma} = 6, \sigma_1 = 1, \sigma_2 = 3, \sigma_3 = 2$ のとき:          |   |   |
| $\hat{p}_{21} = [\underline{p}_{21}, \bar{p}_{21}]$                              | $\hat{p}_{22} = [\underline{p}_{22}, \bar{p}_{22}]$ | $\hat{p}_{23} = [\underline{p}_{23}, \bar{p}_{23}]$ |
| $[0.187, 0.260]$   | $[0.400, 0.489]$                                    | $[0.292, 0.376]$                                    |
| 状態 3, $\hat{\sigma} = 7, \sigma_1 = 1, \sigma_2 = 2, \sigma_3 = 4$ のとき:          |   |   |
| $\hat{p}_{31} = [\underline{p}_{31}, \bar{p}_{31}]$                              | $\hat{p}_{32} = [\underline{p}_{32}, \bar{p}_{32}]$ | $\hat{p}_{33} = [\underline{p}_{33}, \bar{p}_{33}]$ |
| $[0.168, 0.235]$   | $[0.262, 0.334]$                                    | $[0.458, 0.542]$                                    |

$\mathcal{Q} = \langle Q, \bar{Q} \rangle$  の第  $i$  行目に関する凸多面体を  $\hat{q}_i$  ( $i = 1, 2, 3$ ) とおくとき, その端点の集合  $\text{ext}(\hat{q}_i)$  はそれぞれ以下のようになる:  $\text{ext}(\hat{q}_1) = \{(0.437, 0.187, 0.376), (0.4, 0.224, 0.376), (0.448, 0.26, 0.292), (0.489, 0.219, 0.292), (0.4, 0.26, 0.34), (0.489, 0.187, 0.324)\}$ ,  $\text{ext}(\hat{q}_2) = \{(0.187, 0.437, 0.376), (0.224, 0.4, 0.376), (0.26, 0.448, 0.292), (0.219, 0.489, 0.292), (0.26, 0.4, 0.34), (0.187, 0.489, 0.324)\}$ ,  $\text{ext}(\hat{q}_3) = \{(0.196, 0.262, 0.542), (0.168, 0.29, 0.542), (0.208, 0.334, 0.458), (0.235, 0.307, 0.458), (0.168, 0.334, 0.498), (0.235, 0.262, 0.503)\}$ を得る。

$\beta = 0.9, \mathbf{r} = (3, 1, 2)', F \ni f(\text{固定})$  として  $\underline{L}(f)\mathbf{x} = \mathbf{r}(f) + \beta \min_{Q \in \mathcal{Q}(f)} Q\mathbf{x}, \bar{L}(f)\mathbf{x} = \mathbf{r}(f) + \beta \max_{Q \in \mathcal{Q}(f)} Q\mathbf{x}$  の不動点を求めてみると,  $\underline{\phi}(f) = (20.003, 17.508, 18.643), \bar{\phi}(f) = (21.732, 19.232, 20.339)$ を得る。従って, Theorem 3.1 から  $\phi(f|\mathcal{Q}(f)) = [\underline{\phi}(f), \bar{\phi}(f)] = \phi(f|\mathcal{Q}(f)) = ([20.003, 21.732], [17.508, 19.232], [18.643, 20, 339])$ を得る。真の推移確率行列を  $Q = (q_{1.}, q_{2.}, q_{3.})', q_{1.} = (\frac{1}{2}, \frac{1}{6}, \frac{1}{3}), q_{2.} = (\frac{1}{3}, \frac{1}{3}, \frac{1}{3}), q_{3.} = (\frac{2}{5}, \frac{2}{5}, \frac{1}{5})$  であるときの value function の値は  $\phi = (22.469, 20.116, 21.135)$  である。

## 7 区間ベイズ MDPs

最初に, 区間推定 MDPs  $\{\mathcal{Q}\}$  の  $\mathcal{Q} \in \mathcal{M}_n$  に関する連続性を証明する。次に, 事前情報を区間ベイズ法によって処理したデータを使って区間ベイズ MDPs を定義する。

まず,  $\mathcal{Q} = \langle Q, \bar{Q} \rangle \in \mathcal{M}_n$  の  $Q, \bar{Q} \in \mathbb{R}_+^{n \times n}$  の連続性について示そう。次が成り立つ。ただし, 収束は各空間に対応してユークリッド距離とハウスドルフ距離に対応している。

**Lemma 7.1.** (i)  $\underline{Q}_t \downarrow \underline{Q}, \bar{Q}_t \uparrow \bar{Q}$  ( $t \rightarrow \infty$ ),  $\langle \underline{Q}_t, \bar{Q}_t \rangle \neq \emptyset$  ( $t \geq 1$ ) とするとき,  $\langle \underline{Q}_t, \bar{Q}_t \rangle \xrightarrow{\rho} \langle \underline{Q}, \bar{Q} \rangle$  ( $t \rightarrow \infty$ ).  
(ii)  $\underline{Q}_t \uparrow \underline{Q}, \bar{Q}_t \downarrow \bar{Q}$  ( $t \rightarrow \infty$ ),  $\langle \underline{Q}, \bar{Q} \rangle \neq \emptyset$  ( $t \geq 1$ ) とするとき,  $\langle \underline{Q}_t, \bar{Q}_t \rangle \xrightarrow{\rho} \langle \underline{Q}, \bar{Q} \rangle$  ( $t \rightarrow \infty$ ).

*Proof.* (i)  $\langle \underline{Q}_t, \bar{Q}_t \rangle \uparrow$  より  $\{\langle \underline{Q}_t, \bar{Q}_t \rangle\}$  は収束 (i.e.,  $\limsup_{t \rightarrow \infty} \langle \underline{Q}_t, \bar{Q}_t \rangle = \liminf_{t \rightarrow \infty} \langle \underline{Q}_t, \bar{Q}_t \rangle$ ) して,  $\lim_{t \rightarrow \infty} \langle \underline{Q}_t, \bar{Q}_t \rangle = \overline{\cup_{t=1}^{\infty} \langle \underline{Q}_t, \bar{Q}_t \rangle}$  (閉包) であって,  $\langle \underline{Q}_t, \bar{Q}_t \rangle \subset \langle \underline{Q}, \bar{Q} \rangle$  for all  $t \geq 1$  より

$$(7.1) \quad \lim_{t \rightarrow \infty} \langle \underline{Q}_t, \bar{Q}_t \rangle \subset \langle \underline{Q}, \bar{Q} \rangle$$

である.  $\forall Q \in \langle \underline{Q}, \bar{Q} \rangle$  を取る.  $\langle \underline{Q}_t, \bar{Q}_t \rangle \neq \emptyset$  ( $t \geq 1$ ) より  $\exists Q' \in \langle \underline{Q}_t, \bar{Q}_t \rangle$  ( $t \geq 1$ ) であって  $\langle \underline{Q}_t, \bar{Q}_t \rangle \subset \langle \underline{Q}, \bar{Q} \rangle$  より,  $Q' \in \langle \underline{Q}, \bar{Q} \rangle$ .  $0 \leq \alpha \leq 1$  に対して,  $Q(\alpha) = \alpha Q + (1 - \alpha)Q'$  とする. このとき,  $Q(0) = Q' \in \langle \underline{Q}, \bar{Q} \rangle, Q(1) = Q \in \langle \underline{Q}, \bar{Q} \rangle$  かつ Lemma 2.1(i) より  $\langle \underline{Q}, \bar{Q} \rangle$  は凸集合であるから

$$(7.2) \quad Q(\alpha) \in \langle \underline{Q}, \bar{Q} \rangle \quad \text{for all } 0 \leq \alpha \leq 1.$$

$Q(0) = Q' \in \langle \underline{Q}_t, \bar{Q}_t \rangle$  for all  $t \geq 1$  に注意して,  $\alpha_t := \sup\{\alpha | Q(\alpha) \in \langle \underline{Q}_t, \bar{Q}_t \rangle\}$  ( $t \geq 1$ ) とすると,  $\langle \underline{Q}_t, \bar{Q}_t \rangle$  はコンパクトな凸集合であり, かつ  $\langle \underline{Q}_t, \bar{Q}_t \rangle \subset \langle \underline{Q}_{t+1}, \bar{Q}_{t+1} \rangle$  ( $t \geq 1$ ) より,  $Q(\alpha) \in \langle \underline{Q}_t, \bar{Q}_t \rangle$  for all  $0 \leq \alpha \leq \alpha_t$  と  $\alpha_t \leq \alpha_{t+1} \leq 1$  ( $t \geq 1$ ) が成り立つ.

$\alpha^* := \lim_{t \rightarrow \infty} \alpha_t$  とする.  $\alpha^* < 1$  とすると  $\alpha^* < \alpha \leq 1$  なる  $\alpha$  に対して  $Q(\alpha) \notin \langle \underline{Q}, \bar{Q} \rangle$  となるが, これは式 (7.2) に矛盾する. 故に,  $\alpha^* = 1$ . 従って,  $Q(\alpha_t) \in \langle \underline{Q}_t, \bar{Q}_t \rangle \rightarrow Q(1) = Q$  ( $t \rightarrow \infty$ ). これは,  $\limsup_{t \rightarrow \infty} \langle \underline{Q}_t, \bar{Q}_t \rangle \supset \langle \underline{Q}, \bar{Q} \rangle$  を意味する.  $\limsup_{t \rightarrow \infty} \langle \underline{Q}_t, \bar{Q}_t \rangle = \lim_{t \rightarrow \infty} \langle \underline{Q}_t, \bar{Q}_t \rangle$  であるから

$$(7.3) \quad \lim_{t \rightarrow \infty} \langle \underline{Q}_t, \bar{Q}_t \rangle \supset \langle \underline{Q}, \bar{Q} \rangle.$$

式 (7.1) と式 (7.3) より  $\lim_{t \rightarrow \infty} \langle \underline{Q}_t, \bar{Q}_t \rangle = \langle \underline{Q}, \bar{Q} \rangle$ . これで (i) が示された. (ii) は明らかに成り立つ. ■

上の Lemma 7.1 を用いて次が示される.

**Theorem 7.1.**  $\underline{Q}_t \rightarrow \underline{Q}, \bar{Q}_t \rightarrow \bar{Q}$  ( $t \rightarrow \infty$ ),  $\mathcal{Q}_t := \langle \underline{Q}_t, \bar{Q}_t \rangle \neq \emptyset$  ( $t \geq 1$ ),  $\mathcal{Q} := \langle \underline{Q}, \bar{Q} \rangle$  とする. このとき, 次が成り立つ: (i)  $\mathcal{Q}_t \rightarrow \mathcal{Q}$  ( $t \rightarrow \infty$ ). (ii)  $\phi(f|\mathcal{Q}_t) \rightarrow \phi(f|\mathcal{Q})$  ( $t \rightarrow \infty$ ) ( $f \in F$ ).

*Proof.* (i)  $\forall \varepsilon > 0$  に対して  $\limsup_{t \rightarrow \infty} \mathcal{Q}_t = \limsup_{t \rightarrow \infty} \langle \underline{Q}_t, \bar{Q}_t \rangle \subset \langle \underline{Q} - \varepsilon E, \bar{Q} + \varepsilon E \rangle$ , ただし,  $E = (e_{ij})$  は  $e_{ij} = 1$  ( $1 \leq i, j \leq n$ ) とする. このとき,  $\liminf_{t \rightarrow \infty} \mathcal{Q}_t = \liminf_{t \rightarrow \infty} \langle \underline{Q}_t, \bar{Q}_t \rangle \supset \langle \underline{Q} + \varepsilon E, \bar{Q} - \varepsilon E \rangle$  であつて, ここで  $\varepsilon \rightarrow 0$  とすると Lemma 7.1 より  $\lim_{\varepsilon \rightarrow 0} \langle \underline{Q} - \varepsilon E, \bar{Q} + \varepsilon E \rangle = \lim_{\varepsilon \rightarrow 0} \langle \underline{Q} + \varepsilon E, \bar{Q} - \varepsilon E \rangle = \langle \underline{Q}, \bar{Q} \rangle$ . 故に,  $\{\mathcal{Q}_t\}$  は収束して  $\lim_{t \rightarrow \infty} \mathcal{Q}_t = \mathcal{Q}$  を得る. (ii) (i) より,  $\mathcal{Q}_t \rightarrow \mathcal{Q}$  だから,  $\mathbf{x} \in \mathbb{R}_+$  について  $\min_{Q \in \mathcal{Q}_t} Q\mathbf{x} \rightarrow \min_{Q \in \mathcal{Q}} Q\mathbf{x}, \max_{Q \in \mathcal{Q}_t} Q\mathbf{x} \rightarrow \max_{Q \in \mathcal{Q}} Q\mathbf{x}$  ( $t \rightarrow \infty$ ). 従って  $\underline{L}_t(f)\underline{\phi}_t(f) \rightarrow \underline{L}(f)\underline{\phi}(f), \bar{L}_t(f)\bar{\phi}_t(f) \rightarrow \bar{L}(f)\bar{\phi}(f)$  ( $t \rightarrow \infty$ ). 定理 3.1 から  $\phi(f|\mathcal{Q}) = [\underline{\phi}(f), \bar{\phi}(f)]$  であるから  $\phi(f|\mathcal{Q}_t) \rightarrow \phi(f|\mathcal{Q})$  ( $t \rightarrow \infty$ ) を得る. ■

真の推移確率行列  $Q \in P(S|S \times A)$  による MDPs $\{Q\}$  の  $t$  期の状態と行動をそれぞれ  $X_t, \Delta_t$  ( $t \geq 0$ ) で表し,  $t$  期までの履歴を  $H_t = (X_0, \Delta_0, X_1, \Delta_1, \dots, X_t)$  とする. 任意の  $i, j \in S, a \in A$  に対して

$$(7.4) \quad N_T(j|i, a, H_T) := \sum_{t=0}^{T-1} I_{\{X_t=i, \Delta_t=a, X_{t+1}=j\}} \quad (T \geq 1)$$

とおく. 各  $i \in S, a \in A$  に対して, 多項分布の生起確率  $\{p_j = p_{ij}(a), (1 \leq j \leq n)\}$  に対する観測値  $\{N_T(j|i, a, H_T), 1 \leq j \leq n\}$  によるベイズ区間を  $\mathcal{Q}(H_T) = \langle \underline{Q}(H_T), \bar{Q}(H_T) \rangle = \langle (q_{ij}(a|H_T)), (\bar{q}_{ij}(a|H_T)) \rangle$  とする. すなわち,  $\underline{Q}(H_T) := (q_{ij}(a|H_T) : i, j \in S, a \in A) \in \mathbb{R}_+^{n \times nk}, \bar{Q}(H_T) := (\bar{q}_{ij}(a|H_T) : i, j \in S, a \in A) \in \mathbb{R}_+^{n \times nk}$  として,  $\mathcal{Q}(H_T) = \langle \underline{Q}(H_T), \bar{Q}(H_T) \rangle$  とする.  $Q \in P(S|S \times A)$  に対して, MDPs $\{Q\}$  を事前情報  $H_T$  の区間ベイズ  $\mathcal{Q}(H_T)$  で推定した MDPs を区間ベイズ MDPs $\{\mathcal{Q}(H_T)\}$  と言う.  $N_T(i, a|H_T) := \sum_{j \in S} N_T(j|i, a, H_T)$  ( $i \in S, a \in A$ ) とおく.

区間ベイズの性質 ([1]) および Theorem 7.1 を用いて次の結果を得る.

**Theorem 7.2.**  $\{X_0, \Delta_0, X_1, \Delta_1, \dots\}$  を  $MDPs\{Q\}$  からの過程とする. 任意の  $i \in S, a \in A$  に対して, 確率 1 で  $N_T(i, a|H_T) \rightarrow \infty$  ( $T \rightarrow \infty$ ) とする. このとき, 確率 1 で区間ベイズ  $MDPs\{\mathcal{Q}(H_T)\}$  は  $MDPs\{Q\}$  に収束する, すなわち, 次が成り立つ. (i)  $\mathcal{Q}(H_T) \rightarrow \{Q\}$  ( $T \rightarrow \infty$ ), (ii)  $\phi(f|\mathcal{Q}(H_T)) \rightarrow \phi(f|Q)$  ( $T \rightarrow \infty$ ), ( $f \in F$ ).

*Proof.* Theorem 7.2 の条件が成り立てば, [1] の定理 5.2 より  $\underline{Q}(H_T), \bar{Q}(H_T) \rightarrow Q$  ( $T \rightarrow \infty$ ) を得る. 従つて, Theorem 7.1 より (i),(ii) が成り立つことがわかる. ■

## References

- [1] L. De Robertis and J. A. Hartigan. Bayesian inference using intervals of measures. *Ann. Statist.*, 9(2):235–244, 1981.
- [2] B. Doshi and S. E. Shreve. Strong consistency of a modified maximum likelihood estimator for controlled Markov chains. *J. Appl. Probab.*, 17(3):726–734, 1980.
- [3] N. Furukawa. Characterization of optimal policies in vector-valued Markovian decision processes. *Math. Oper. Res.*, 5(2):271–279, 1980.
- [4] D. J. Hartfiel. *Markov set-chains*, volume 1695 of *Lecture Notes in Mathematics*. Springer-Verlag, Berlin, 1998.
- [5] O. Hernández-Lerma. *Adaptive Markov control processes*, volume 79 of *Applied Mathematical Sciences*. Springer-Verlag, New York, 1989.
- [6] M. Horiguchi. Examples for Bayesian approach to uncertain MDPs,  
[URL](http://www.math.kanagawa-u.ac.jp/~horiguchi/) <http://www.math.kanagawa-u.ac.jp/~horiguchi/>
- [7] T. Iki, M. Horiguchi, M. Yasuda, and M. Kurano. A learning algorithm for communicating markov decision processes with unknown transition matrices. *Bulletin of Informatics and Cybernetics*, 39:11–24, 2007.
- [8] M. Kurano, J. Song, M. Hosaka, and Y. Huang. Controlled Markov set-chains with discounting. *J. Appl. Probab.*, 35(2):293–302, 1998.
- [9] M. Kurano, M. Yasuda, and J. Nakagami. Interval methods for uncertain Markov decision processes. In *Markov processes and controlled Markov chains (Changsha, 1999)*, pages 223–232. Kluwer, 2002.
- [10] K. Kuratowski. *Topology. Vol. I.* New edition, revised and augmented. Translated from the French by J. Jaworowski. Academic Press, 1966.
- [11] P. Mandl. Estimation and control in Markov chains. *Advances in Appl. Probability*, 6:40–60, 1974.
- [12] J. J. Martin. *Bayesian decision problems and Markov chains*. Publications in Operations Research, No. 13. John Wiley & Sons Inc., 1967.
- [13] M. L. Puterman. *Markov decision processes: discrete stochastic dynamic programming*. John Wiley & Sons Inc., 1994.
- [14] E. Solan. Continuity of the value of competitive Markov decision processes. *J. Theoret. Probab.*, 16(4):831–845 (2004), 2003.
- [15] K. M. van Hee. *Bayesian control of Markov chains*, volume 95 of *Mathematical Centre Tracts*. Mathematisch Centrum, 1978.
- [16] S. S. Wilks. *Mathematical statistics*. A Wiley Publication in Mathematical Statistics. John Wiley & Sons Inc., 1962. 田中英之, 岩本誠一 (訳), 「数理統計学・増訂新版 1,2」, 1971, 1972 年, 東京図書.