

# 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・理学研究科 安田正實 (Masami YASUDA)

Faculty of Science, Chiba University  
蕨野正美 (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, 6, 10]) とベイズ推定法を用いる場合 (cf. [14, 5, 11]) 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くつかの数値例を与える。Kurano et al.[7, 8] 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内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

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

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:

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

$n \times n$  型の確率行列の区間集合全体を  $\mathcal{M}_n$  で表す。

$$(2.3) \quad \mathcal{M}_n = \{ \langle \underline{Q}, \bar{Q} \rangle \mid \langle \underline{Q}, \bar{Q} \rangle \neq \emptyset, \underline{Q} \preceq \bar{Q}, \underline{Q}, \bar{Q} \in \mathbb{R}_+^{n \times n} \}$$

$\mathcal{M}_n \ni \mathcal{Q}_1, \mathcal{Q}_2$  に対する積  $\mathcal{Q}_1, \mathcal{Q}_2$  を次で定める。

$$(2.4) \quad \mathcal{Q}_1 \mathcal{Q}_2 = \{ Q_1 Q_2 \mid Q_1 \in \mathcal{Q}_1, Q_2 \in \mathcal{Q}_2 \}$$

また,  $\mathcal{Q} \in \mathcal{M}_n$  に対する多重積は逐次的に定義される:

$$(2.5) \quad \mathcal{Q}^k = \mathcal{Q}^{k-1}\mathcal{Q} \ (k \geq 2).$$

$\mathbb{C}(\mathbb{R}_+)$  を  $\mathbb{R}_+$  の有界閉区間の全体とする. また,  $\mathbb{C}(\mathbb{R}_+)^n$  を  $\mathbb{C}(\mathbb{R}_+)$  の要素を成分に持つ  $n$  次元列ベクトルの全体とする:

$$(2.6) \quad \mathbb{C}(\mathbb{R}_+)^n = \{D = (D_1, D_2, \dots, D_n)' \mid D_i \in \mathbb{C}(\mathbb{R}_+) \ (1 \leq i \leq n)\}$$

ただし,  $\mathbf{d}'$  はベクトル  $\mathbf{d}$  の転置を表す.

$\mathbb{C}(\mathbb{R}_+)^n$  上の算法 (加法, スカラ一倍) は次で定める:  $D = (D_1, D_2, \dots, D_n)', E = (E_1, E_2, \dots, E_n)' \in \mathbb{C}(\mathbb{R}_+)^n, h \in \mathbb{R}_+^n, \lambda \in \mathbb{R}_+$  に対して,

$$(2.7) \quad D + E = \{d + e \mid d \in D, e \in E\}, h + D = \{h + d \mid d \in D\}, \lambda D = \{\lambda d \mid d \in D\}.$$

$D = ([\underline{d}_1, \bar{d}_1], [\underline{d}_2, \bar{d}_2], \dots, [\underline{d}_n, \bar{d}_n])' \in \mathbb{C}(\mathbb{R}_+)^n$  を  $D = [\underline{d}, \bar{d}]$  と記す. ただし,  $\underline{d} = (\underline{d}_1, \underline{d}_2, \dots, \underline{d}_n) \in \mathbb{R}_+^n, \bar{d} = (\bar{d}_1, \bar{d}_2, \dots, \bar{d}_n) \in \mathbb{R}_+^n$  とする.  $D = (D_1, D_2, \dots, D_n)' \in \mathbb{C}(\mathbb{R}_+)^n$  と部分集合  $G \subset \mathbb{R}_+^{1 \times n}$  に対して, その積  $GD$  を次で定める:

$$(2.8) \quad GD = \{gd \mid g = (g_1, g_2, \dots, g_n) \in G, d = (d_1, d_2, \dots, d_n)' \in D, d_i \in D_i \ (1 \leq i \leq n)\}$$

次が成り立つ.

**Lemma 2.1.** ([4, 7])

- (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  $D = (D_1, D_2, \dots, D_n) \in \mathbb{C}(\mathbb{R}_+)^n$  に対して  $GD \in \mathbb{C}(\mathbb{R}_+)$  である.

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

$$(2.9) \quad \begin{cases} [c_1, c_2] \preceq [d_1, d_2] & (c_i \leq d_i \ (i = 1, 2) \text{ のとき}) \\ [c_1, c_2] \prec [d_1, d_2] & ([c_1, c_2] \preceq [d_1, d_2] \text{ かつ } [c_1, c_2] \neq [d_1, d_2] \text{ のとき}) \end{cases}$$

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

$$(2.10) \quad \begin{cases} \mathbf{v} \preceq \mathbf{w} & (v_i \leq w_i \ (1 \leq i \leq n) \text{ のとき}) \\ \mathbf{v} \prec \mathbf{w} & (\mathbf{v} \prec \mathbf{w} \text{ かつ } \mathbf{v} \neq \mathbf{w} \text{ のとき}) \end{cases}$$

$\mathbb{R}_+^n$  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:

$$(2.11) \quad \rho(D_1, D_2) = \max\{\sup_{x \in D_1} \inf_{y \in D_2} \|x - y\|, \sup_{y \in D_2} \inf_{x \in D_1} \|x - y\|\}.$$

ただし,  $\|\cdot\|$  は  $\mathbb{R}^n$  におけるヨ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:

$$(2.12) \quad P(S) := \{p = (p_1, p_2, \dots, p_n) \in \mathbb{R}_+^n \mid \sum_{i \in S} p_i = 1\},$$

$$(2.13) \quad P(S|S) := \{q = (q_{ij} : i, j \in S) \in \mathbb{R}_+^{n \times n} \mid \sum_{j \in S} q_{ij} = 1 \ (i \in S)\},$$

$$(2.14) \quad P(S|S \times A) := \{Q = (q_{ij}(a) : i, j \in S, a \in A) \in \mathbb{R}_+^{kn \times n} \mid q_{i.}(a) \in P(s) \ (i \in S, a \in A)\}.$$

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  $\{S, A, Q, \mathbf{r}\}$  を考え (cf. [12]), 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:

$$(2.15) \quad \phi(f|Q) = \sum_{t=0}^{\infty} (\beta Q(f))^t \mathbf{r}(f),$$

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:

$$(2.16) \quad L(f)\mathbf{x} = \mathbf{r}(f) + \beta Q(f)\mathbf{x}, \quad \mathbf{x} = (x_1, x_2, \dots, x_n)' \in \mathbb{R}_+^n.$$

このとき, 次の基本補題が知られている.

**Lemma 2.2.** (cf. [12])

(i)  $L(f)$  は単調増加および縮小写像である. すなわち,

$$\begin{aligned} \mathbf{x} \leqq \mathbf{x}' \text{ ならば } L(f)\mathbf{x} \leqq L(f)\mathbf{x}' \text{ (componentwise),} \\ \|L(f)\mathbf{x} - L(f)\mathbf{x}'\| \leqq \beta \|\mathbf{x} - \mathbf{x}'\| \quad (\mathbf{x}, \mathbf{x}' \in \mathbb{R}_+^n), \end{aligned}$$

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

$$L(f)^t \mathbf{x} \rightarrow \phi(f|Q) \quad (t \rightarrow \infty)$$

が成り立つ.

### 3 区間推定 MDPs とパレート最適

本節では, MDP( $S, A, Q, r$ ) 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,

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

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

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 (2.15) で与えられている.

ここで,  $\phi(f|\mathcal{Q}) \in \mathbb{C}(\mathbb{R}_+)^n$  であることを示そう.  $\mathcal{L}$  を  $\mathbb{C}(\mathbb{R}_+)^n$  から  $\mathbb{C}(\mathbb{R}_+)^n$  への写像で次のように定める:

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

ただし, 式 (3.4)において  $Q(f) = \langle \underline{Q}(f), \bar{Q}(f) \rangle$ ,  $\underline{Q}(f) = (\underline{q}_{ij}(f(i))) \in \mathbb{R}_+^{n \times n}$ ,  $\bar{Q}(f) = (\bar{q}_{ij}(f(i))) \in \mathbb{R}_+^{n \times n}$  である.

Lemma 2.1 により  $\mathcal{L}(f)\mathbf{v} \in \mathbb{C}(\mathbb{R}_+)^n$  ( $\mathbf{v} \in \mathbb{C}(\mathbb{R}_+)^n$ ) 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

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

$$(3.6) \quad \bar{L}(f)\mathbf{x} = \mathbf{r}(f) + \beta \max_{Q \in \mathcal{Q}(f)} Q\mathbf{x}.$$

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.

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)  $\phi(f|\mathcal{Q}) \in \mathbb{C}(\mathbb{R}_+)^n$  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  $\mathbf{v} \in \mathbb{C}(\mathbb{R}_+)^n$  に対して

$$\mathcal{L}(f)^\ell \mathbf{v} \rightarrow \phi(f|\mathcal{Q}) (\ell \rightarrow \infty)$$

(ii)  $\phi(f|\mathcal{Q}) = [\underline{\phi}(f), \bar{\phi}(f)]$  とするとき,  $\underline{\phi}(f), \bar{\phi}(f)$  はそれぞれ  $\underline{L}(f), \bar{L}(f)$  の唯一の不動点である.

$f^* \in F$  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  $\phi(f|\mathcal{Q}) \prec \phi(g|\mathcal{Q})$ .

$D \subset \mathbb{C}(\mathbb{R}_+)^n$  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.  $\mathbf{u} \in \mathbb{C}(\mathbb{R}_+)^n$  に対して次を定める:

$$(3.7) \quad \mathcal{L}(\mathbf{u}) := (\mathcal{L}(\mathbf{u})_1, \mathcal{L}(\mathbf{u})_2, \dots, \mathcal{L}(\mathbf{u})_n)',$$

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

$$(3.8) \quad \mathbf{u} \in \mathcal{L}(\mathbf{u}), \mathbf{u} \in \mathbb{C}(\mathbb{R}_+)^n$$

## 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。

はじめに、ガンマ関数  $\Gamma(x)$  ( $x > 0$ ) とベータ関数  $B(x, y)$  ( $x, y > 0$ ) と簡単な性質についてまとめておく。

$$\text{ガンマ関数: } \Gamma(x) = \int_0^\infty t^{x-1} e^{-t} dt \quad (x > 0), \quad \text{ベータ関数: } B(x, y) = \int_0^1 t^{x-1} (1-t)^{y-1} dt \quad (x, y > 0)$$

性質:

- $\Gamma(x+1) = x\Gamma(x)$ ,  $\Gamma(1) = 1$ ,  $\Gamma\left(\frac{1}{2}\right) = \sqrt{\pi}$ ,  $B(x, y) = \frac{\Gamma(x)\Gamma(y)}{\Gamma(x+y)}$

ディリクレ分布 ([15]):

ベータ分布  $\tilde{B}(\nu_1, \nu_2)$  の p.d.f.

$$(4.1) \quad f(x) = \frac{\Gamma(\nu_1 + \nu_2)}{\Gamma(\nu_1)\Gamma(\nu_2)} x^{\nu_1-1} (1-x)^{\nu_2-1}$$

を拡張して、 $k$ -変数ディリクレ分布 ( $k$ -variable Dirichlet distribution)  $\tilde{D}(\nu_1, \dots, \nu_k; \nu_{k+1})$  の p.d.f. を次のように定義する:

$$(4.2) \quad f(x_1, \dots, x_k) = \frac{\Gamma(\nu_1 + \dots + \nu_{k+1})}{\Gamma(\nu_1) \dots \Gamma(\nu_{k+1})} x_1^{\nu_1-1} \dots x_k^{\nu_k-1} (1 - x_1 - x_2 - \dots - x_k)^{\nu_{k+1}-1}$$

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 とする。 $\nu_i \in \mathbb{R}$  は  $\nu_i > 0$  とする。

$$(4.3) \quad \tilde{D}(\nu_1, \dots, \nu_k; \nu_{k+1}) = \int \dots \int_{S_k} f(x_1, \dots, x_k) dx_1 \dots dx_k$$

と表すとき次が成り立つ:

- $k = 1$  のとき、 $\tilde{D}(\nu_1; \nu_2)$  は  $\tilde{B}(\nu_1, \nu_2)$  である。(定義から明らか)
- $k \geq 2$  のとき、式 (4.2) は確かに確率密度関数を表し、ベータ関数(ベータ分布でないことに注意)の値  $B(x, y)$  を用いて

$$(4.4) \quad \begin{aligned} & \tilde{D}(\nu_1, \nu_2, \dots, \nu_k; \nu_{k+1}) \\ &= \frac{\Gamma(\nu_1 + \dots + \nu_{k+1})}{\Gamma(\nu_1)\Gamma(\nu_2) \dots \Gamma(\nu_{k+1})} B(\nu_1, \nu_2 + \dots + \nu_{k+1}) B(\nu_2, \nu_3 + \dots + \nu_{k+1}) \dots B(\nu_k, \nu_{k+1}) \end{aligned}$$

が成り立つ。

実際、確率密度関数となることは以下のようない変数変換をとることで示される:

$$(4.5) \quad \begin{aligned} x_1 &= \theta_1, x_2 = \theta_2(1-x_1) = \theta_2(1-\theta_1), x_3 = \theta_3(1-x_1-x_2) = \theta_3(1-\theta_1)(1-\theta_2), \\ &\vdots \\ x_k &= \theta_k(1-x_1-x_2-\dots-x_{k-1}) = \theta_k(1-\theta_1)(1-\theta_2)\dots(1-\theta_{k-1}) \end{aligned}$$

とすると  $S_k = \{(x_1, \dots, x_k) : x_1 + x_2 + \dots + x_k \leq 1, x_i \geq 0\}$  は  $k$ -次元直方体  $U_k := \{(\theta_1, \dots, \theta_k) : 0 \leq \theta_i \leq 1, i = 1, 2, \dots, k\}$  に 1 対 1 対応で移される。ヤコビアンは

$$\left| \frac{\partial(x_1, \dots, x_k)}{\partial(\theta_1, \dots, \theta_k)} \right| = (1 - \theta_1)^{k-1} (1 - \theta_2)^{k-2} \cdots (1 - \theta_{k-2})^2 (1 - \theta_{k-1})$$

となる。従って、

$$\begin{aligned}
 \tilde{D}(\nu_1, \dots, \nu_k; \nu_{k+1}) &= \int \cdots \int_{S_k} f(x_1, \dots, x_k) dx_1 \cdots dx_k \\
 &= \int \cdots \int_{U_k} \frac{\Gamma(\nu_1 + \nu_2 + \cdots + \nu_{k+1})}{\Gamma(\nu_1)\Gamma(\nu_2)\cdots\Gamma(\nu_{k+1})} \theta_1^{\nu_1-1} (1 - \theta_1)^{\nu_2+\cdots+\nu_{k+1}-1} \theta_2^{\nu_2-1} (1 - \theta_2)^{\nu_3+\cdots+\nu_{k+1}-1} \cdots \\
 &\quad \cdots \theta_k^{\nu_k-1} (1 - \theta_k)^{\nu_{k+1}-1} d\theta_1 d\theta_2 \cdots d\theta_k \\
 (4.6) \quad &= \frac{\Gamma(\nu_1 + \cdots + \nu_{k+1})}{\Gamma(\nu_1)\Gamma(\nu_2)\cdots\Gamma(\nu_{k+1})} B(\nu_1, \nu_2 + \cdots + \nu_{k+1}) B(\nu_2, \nu_3 + \cdots + \nu_{k+1}) B(\nu_k, \nu_{k+1}) \\
 &= \frac{\Gamma(\nu_1 + \cdots + \nu_{k+1})}{\Gamma(\nu_1)\Gamma(\nu_2)\cdots\Gamma(\nu_{k+1})} \frac{\Gamma(\nu_1)\Gamma(\nu_2 + \cdots + \nu_{k+1})}{\Gamma(\nu_1 + \cdots + \nu_{k+1})} \frac{\Gamma(\nu_2)\Gamma(\nu_3 + \cdots + \nu_{k+1})}{\Gamma(\nu_2 + \cdots + \nu_{k+1})} \cdots \frac{\Gamma(\nu_k)\Gamma(\nu_{k+1})}{\Gamma(\nu_k + \nu_{k+1})} \\
 &= 1
 \end{aligned}$$

を得る。

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

$$\begin{aligned}
 D(\nu_1, \nu_2, \dots, \nu_k; \nu_{k+1}) &:= \int \cdots \int_{S_k} x_1^{\nu_1-1} x_2^{\nu_2-1} \cdots x_k^{\nu_k-1} (1 - x_1 - x_2 - \cdots - x_k)^{\nu_{k+1}-1} dx_1 dx_2 \cdots dx_k \\
 (4.7) \quad &= \frac{\Gamma(\nu_1) \cdots \Gamma(\nu_{k+1})}{\Gamma(\nu_1 + \cdots + \nu_{k+1})}
 \end{aligned}$$

であるから、

$$\begin{aligned}
 D(\nu_1, \nu_2, \dots, \nu_k; \nu_{k+1}) &= \frac{\Gamma(\nu_1)\Gamma(\nu_2 + \cdots + \nu_{k+1})}{\Gamma(\nu_1 + \cdots + \nu_{k+1})} \frac{\Gamma(\nu_2)\Gamma(\nu_3) \cdots \Gamma(\nu_{k+1})}{\Gamma(\nu_2 + \cdots + \nu_{k+1})} \\
 &= B(\nu_1, \nu_2 + \cdots + \nu_{k+1}) D(\nu_2, \nu_3, \dots, \nu_k; \nu_{k+1}) = \cdots \\
 (4.8) \quad &= B(\nu_1, \nu_2 + \cdots + \nu_{k+1}) B(\nu_2, \nu_3 + \cdots + \nu_{k+1}) \cdots \\
 &\quad \cdots B(\nu_{k-1}, \nu_k + \cdots + \nu_{k+1}) D(\nu_k; \nu_{k+1}) \\
 &= \prod_{n=1}^{n=k} B\left(\nu_n, \sum_{l=n+1}^{k+1} \nu_l\right)
 \end{aligned}$$

を得る。

$0 < \lambda < 1$  に対して、

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

とする。特に

$$B(\alpha, \beta | \lambda) = D(\alpha; \beta | \lambda) \quad (\alpha, \beta > 0)$$

と表すことにする。式(4.5)と同様の変数変換により

$$(4.10) \quad D(\nu_1, \dots, \nu_k; \nu_{k+1} | \lambda)$$

$$(4.11) \quad = \int_0^\lambda \theta_1^{\nu_1-1} (1-\theta_1)^{\nu_2+\dots+\nu_{k+1}-1} d\theta_1 \int_0^1 \theta_2^{\nu_2-1} (1-\theta_2)^{\nu_3+\dots+\nu_{k+1}-1} d\theta_2 \dots \\ \dots \int_0^1 \theta_k^{\nu_k-1} (1-\theta_k)^{\nu_{k+1}-1} d\theta_k$$

$$(4.12) \quad = B(\nu_1, \nu_2 + \dots + \nu_{k+1} | \lambda) B(\nu_2, \nu_3 + \dots + \nu_{k+1}) B(\nu_3, \nu_4 + \dots + \nu_{k+1}) \dots B(\nu_k, \nu_{k+1})$$

であることがわかる。ここで、 $m, n$ を正の整数とするとき

$$B(m, n | \lambda) = \int_0^\lambda x^{m-1} (1-x)^{n-1} dx \quad (m, n > 0)$$

を  $x = \lambda\theta$  として置換積分してみると

$$(4.13) \quad B(m, n | \lambda) = \int_0^1 (\lambda\theta)^{m-1} (1-\lambda\theta)^{n-1} \lambda d\theta$$

$$(4.14) \quad = \lambda^m \int_0^1 \theta^{m-1} \left( \sum_{i=0}^{n-1} \binom{n-1}{i} (-\lambda\theta)^i \right) d\theta$$

$$(4.15) \quad = \lambda^m \sum_{i=0}^{n-1} \binom{n-1}{i} (-\lambda)^i \int_0^1 \theta^{m+i-1} d\theta = \sum_{i=0}^{n-1} \binom{n-1}{i} (-1)^i \lambda^{m+i} \frac{1}{m+i}.$$

また、

$$(4.16) \quad \frac{d}{d\lambda} B(m, n | \lambda) = \lambda^{m-1} (1-\lambda)^{n-1}$$

であることも注意しておこう。

## 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 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。(※前節の  $k$  次元多面体  $S_k$  との関係は  $P_n$  の  $0 \leq p_n \leq 1$  に関する切片の空間  $\{(p_1, \dots, p_{n-1}) | \sum_{i=1}^{n-1} p_i \leq 1, p_i \geq 0, i = 1, 2, \dots, n-1\}$  は  $S_{n-1}$  に等しい。)

$L(\cdot)$  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 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  $\bar{\sigma} := \sum_{k=1}^n \sigma_k$  回の独立試行実験でそれぞれ state  $i$  が  $\sigma_i$  回起きたことを表す。state  $i$  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  $\sigma$  の p.d.f. は多項分布で表されて

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

となる。

データ  $\sigma$  における事前測度区間を  $[L_\sigma, U_\sigma] = [L_\sigma, kL_\sigma]$  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

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

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。

$U_\sigma = kL_\sigma$  であるから、(5.2), (5.3) は

$$kL_\sigma(p_1 - \underline{\lambda})^- + L_\sigma(p_1 - \underline{\lambda})^+ = 0, kL_\sigma(p_1 - \bar{\lambda})^+ + L_\sigma(p_1 - \bar{\lambda})^- = 0$$

となる。従って、

$$(5.4) \quad k \int_{0 \leq p_1 \leq 1, p \in P_n} \cdots \int (p_1 - \underline{\lambda})^- L_\sigma(dp) + \int_{0 \leq p_1 \leq 1, p \in P_n} \cdots \int (p_1 - \underline{\lambda})^+ L_\sigma(dp) = 0$$

$$(5.5) \quad k \int_{0 \leq p_1 \leq 1, p \in P_n} \cdots \int (p_1 - \bar{\lambda})^+ L_\sigma(dp) + \int_{0 \leq p_1 \leq 1, p \in P_n} \cdots \int (p_1 - \bar{\lambda})^- L_\sigma(dp) = 0$$

であって、

$$(5.6) \quad (p_1 - \lambda)^- = \begin{cases} 0, & (\lambda \leq p_1 \leq 1) \\ p_1 - \lambda, & (0 \leq p_1 < \lambda) \end{cases}, (p_1 - \lambda)^+ = \begin{cases} p_1 - \lambda, & (\lambda < p_1 \leq 1) \\ 0, & (0 \leq p_1 \leq \lambda) \end{cases}$$

に注意すれば、lower bound  $\underline{\lambda}$  と upper bound  $\bar{\lambda}$  に関する  $\lambda$  の方程式は次のようになる：

lower bound  $\underline{\lambda}$

$$(5.7) \quad k \int_{0 \leq p_1 \leq \underline{\lambda}, p \in P_n} \cdots \int (p_1 - \lambda) p_1^{\sigma_1} \cdots p_n^{\sigma_n} dp + \int_{\lambda \leq p_1 \leq 1, p \in P_n} \cdots \int (p_1 - \lambda) p_1^{\sigma_1} \cdots p_n^{\sigma_n} dp = 0$$

upper bound  $\bar{\lambda}$

$$(5.8) \quad k \int_{\lambda \leq p_1 \leq 1, p \in P_n} \cdots \int (p_1 - \lambda) p_1^{\sigma_1} \cdots p_n^{\sigma_n} dp + \int_{0 \leq p_1 \leq \lambda, p \in P_n} \cdots \int (p_1 - \lambda) p_1^{\sigma_1} \cdots p_n^{\sigma_n} dp = 0$$

ここで、前節の結果

$$(5.9) \quad D(\nu_1, \nu_2, \dots, \nu_k; \nu_{k+1}) = B(\nu_1, \nu_2 + \dots + \nu_{k+1}) D(\nu_2, \nu_3, \dots, \nu_k; \nu_{k+1})$$

$$(5.10) \quad D(\nu_1, \nu_2, \dots, \nu_k; \nu_{k+1} | \lambda) = B(\nu_1, \nu_2 + \dots + \nu_{k+1} | \lambda) D(\nu_2, \nu_3, \dots, \nu_k; \nu_{k+1})$$

とベータ関数の性質を利用して、式 (5.7) と式 (5.8) は次の  $\lambda$  に関する多項式の方程式の解であることが示される。

$\bar{\sigma} = \sigma_1 + \sigma_2 + \dots + \sigma_n, p = \sigma_1 + 1, q = \bar{\sigma} - \sigma_1 + n - 1$  とおくと、式 (5.7) と式 (5.8) は結局

$$(5.11) \quad K(p, q, \lambda) := \left( \frac{p}{p+q} - \lambda \right) B(p, q) + (k-1) (B(p+1, q | \lambda) - \lambda B(p, q | \lambda)) = 0$$

$$(5.12) \quad G(p, q, \lambda) := k \left( \frac{p}{p+q} - \lambda \right) B(p, q) - (k-1) (B(p+1, q | \lambda) - \lambda B(p, q | \lambda)) = 0$$

と表される。

**Theorem 5.1.** データ  $\sigma = (\sigma_1, \sigma_2, \dots, \sigma_n), \bar{\sigma} = \sum_{i=1}^n \sigma_i$  とする。事前測度区間を  $[L_\sigma, kL_\sigma]$  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解である。

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

## 6 $p_1$ の事後測度区間 $[\underline{\lambda}, \bar{\lambda}]$ と多項式

前節で求めた  $p_1$  の事後測度区間  $[\underline{\lambda}, \bar{\lambda}]$  の値  $\underline{\lambda}, \bar{\lambda}$  は、具体的には、それぞれ式 (5.11) と (5.12) から次の  $(\bar{\sigma} + n)$  次多項式の解になっている。

$$(6.1) \quad K(p, q, \lambda) = \left( \frac{p}{p+q} - \lambda \right) B(p, q) + (k-1) \left( \sum_{i=0}^{q-1} \binom{q-1}{i} (-1)^{i+1} \lambda^{(p+1)+i} \left( \frac{1}{(p+1+i)(p+i)} \right) \right) = 0$$

$$(6.2) \quad G(p, q, \lambda) = k \left( \frac{p}{p+q} - \lambda \right) B(p, q) - (k-1) \left( \sum_{i=0}^{q-1} \binom{q-1}{i} (-1)^{i+1} \lambda^{(p+1)+i} \left( \frac{1}{(p+1+i)(p+i)} \right) \right) = 0$$

ただし、 $\bar{\sigma} = \sum_{i=1}^n \sigma_i, p = \sigma_1 + 1, q = \bar{\sigma} - \sigma_1 + n - 1$ 。

$K(p, q, \lambda), G(p, q, \lambda)$  はともに狭義単調関数で、 $K(p, q, \lambda)$  は上に凸、 $G(p, q, \lambda)$  は下に凸である。

$$(6.3) \quad \frac{dK}{d\lambda} = -B(p, q) + (k-1) (\lambda^p (1-\lambda)^{q-1} - B(p, q|\lambda) - \lambda^p (1-\lambda)^{q-1}) \\ = -B(p, q) - (k-1) B(p, q|\lambda) < 0$$

$$(6.4) \quad \frac{d^2K}{d\lambda^2} = -(k-1) \lambda^p (1-\lambda)^{q-1} < 0$$

$$(6.5) \quad \frac{dG}{d\lambda} = -kB(p, q) - (k-1) (\lambda^p (1-\lambda)^{q-1} - B(p, q|\lambda) - \lambda^p (1-\lambda)^{q-1}) \\ = -kB(p, q) + (k-1) B(p, q|\lambda) \\ \leq -kB(p, q) + (k-1) B(p, q) = -B(p, q) < 0$$

$$(6.6) \quad \frac{d^2G}{d\lambda^2} = (k-1) \lambda^{p-1} (1-\lambda)^{q-1} > 0$$

また、 $\bar{\lambda}, \underline{\lambda}$  が一意の解であることは [1] の結果から明らかであるが、

$$(6.7) \quad G(p, q, 0) = kB(p+1, q) > 0, G(p, q, 1) = -\frac{q}{p+q} B(p, q) < 0$$

$$(6.8) \quad K(p, q, 0) = B(p+1, q) > 0, K(p, q, 1) = -\frac{kq}{p+q} B(p, q) < 0$$

と  $G, K$  の単調性から  $\lambda$  に関して  $[0, 1]$  では必ず一つの解のみを持つことがわかる。

## 7 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  $\bar{\sigma} = 6$  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、

$$\bar{\sigma} = \sigma_1 + \sigma_2 + \sigma_3 = 6, p = \sigma_1 + 1 = 4, q = \sigma_2 + \sigma_3 + (n-1) = 5$$

のデータが得られているとする。

式 (6.2) の  $\bar{\lambda}$  に関する多項式は

$$(7.1) \quad 2 \left( \frac{4}{6+3} - \lambda \right) B(4, 5) - \left( \sum_{i=0}^4 \binom{4}{i} (-1)^{i+1} \lambda^{5+i} \left( \frac{1}{(4+i)(5+i)} \right) \right) = 0$$

より

$$(7.2) \quad 8 - 18\lambda + \lambda^5 (126 - 336\lambda + 360\lambda^2 - 180\lambda^3 + 35\lambda^4) = 0$$

となる。このとき、解  $\bar{\lambda} \approx 0.489$  を得る。

また、式(6.1)の $\lambda$ に関する多項式は

$$(7.3) \quad \left( \frac{4}{6+3} - \lambda \right) B(4,5) + \left( \sum_{i=0}^4 \binom{4}{i} (-1)^{i+1} \lambda^{5+i} \left( \frac{1}{(4+i)(5+i)} \right) \right) = 0$$

より

$$(7.4) \quad 4 - 9\lambda - \lambda^5(126 - 336\lambda + 360\lambda^2 - 180\lambda^3 + 35\lambda^4) = 0$$

となる。このとき、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。

同様にして、状態数 $N = 3$ 、実験回数 $\bar{\sigma}$ をいろいろ変化させたときの $p_1, p_2, p_3$ の事後測度区間をまとめると(Table 1)。 $\sigma_1, \sigma_2, \sigma_3$ はそれぞれ状態*i*での観測値(回数)とする。

Table 1: 数値実験例(状態数 $N = 3$ )

$\bar{\sigma} = 6$ (実験回数),  $\sigma_1 = 3, \sigma_2 = 1, \sigma_3 = 2$ のとき,

$p_1 = [\underline{p}_1, \bar{p}_1]$	$p_2 = [\underline{p}_2, \bar{p}_2]$	$p_3 = [\underline{p}_3, \bar{p}_3]$
[0.400, 0.489]	[0.187, 0.260]	[0.292, 0.376]

$\bar{\sigma} = 15, \sigma_1 = 7, \sigma_2 = 3, \sigma_3 = 5$ のとき,

$p_1 = [\underline{p}_1, \bar{p}_1]$	$p_2 = [\underline{p}_2, \bar{p}_2]$	$p_3 = [\underline{p}_3, \bar{p}_3]$
[0.413, 0.476]	[0.197, 0.249]	[0.304, 0.364]

$\bar{\sigma} = 30, \sigma_1 = 16, \sigma_2 = 5, \sigma_3 = 9$ のとき,

$p_1 = [\underline{p}_1, \bar{p}_1]$	$p_2 = [\underline{p}_2, \bar{p}_2]$	$p_3 = [\underline{p}_3, \bar{p}_3]$
[0.491, 0.539]	[0.164, 0.201]	[0.281, 0.325]

$\bar{\sigma} = 50, \sigma_1 = 24, \sigma_2 = 9, \sigma_3 = 17$ のとき,

$p_1 = [\underline{p}_1, \bar{p}_1]$	$p_2 = [\underline{p}_2, \bar{p}_2]$	$p_3 = [\underline{p}_3, \bar{p}_3]$
[0.453, 0.491]	[0.174, 0.204]	[0.322, 0.358]

$k = 1$ 、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

$$p_i = [\underline{p}_i, \bar{p}_i] = \frac{\sigma_i + 1}{\bar{\sigma} + n}$$

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 $p_i$ の周辺分布の期待値に等しい。

ここで、具体的に $L(\cdot)$ :ルベーグ測度、事前測度区間 $[L, kL](k$ は定数)について、 $k = 2$ と考えて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。

状態数 $N = 3, S = \{1, 2, 3\}$ 、policyは固定(deterministic stationary policy)として初期状態 $x_1 = 1$ から状態推移の観測20回で、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

$$\begin{pmatrix} 3 & 1 & 2 \\ 1 & 3 & 2 \\ 1 & 2 & 4 \end{pmatrix}$$

であった。例えば、状態2からの推移では、上の行列の第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 2)。 $\sigma_1, \sigma_2, \sigma_3$ はそれぞれ状態*i*での観測値(推移回数)とする。

Table 2: Intervals of posterior measures

$\bar{\sigma} = 6$ (実験回数),  $\sigma_1 = 3, \sigma_2 = 1, \sigma_3 = 2$ のとき,

$\bar{\sigma} = 6, \sigma_1 = 1, \sigma_2 = 3, \sigma_3 = 2$ のとき,

$\hat{p}_{11} = [\underline{p}_{11}, \bar{p}_{11}]$	$\hat{p}_{12} = [\underline{p}_{12}, \bar{p}_{12}]$	$\hat{p}_{13} = [\underline{p}_{13}, \bar{p}_{13}]$
[0.400, 0.489]	[0.187, 0.260]	[0.292, 0.376]

$\hat{p}_{21} = [\underline{p}_{21}, \bar{p}_{21}]$

$\hat{p}_{22} = [\underline{p}_{22}, \bar{p}_{22}]$

$\hat{p}_{23} = [\underline{p}_{23}, \bar{p}_{23}]$

[0.187, 0.260]

[0.400, 0.489]

[0.292, 0.376]

$\bar{\sigma} = 7, \sigma_1 = 1, \sigma_2 = 2, \sigma_3 = 4$ 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]

$\mathcal{Q} = \langle Q, \bar{Q} \rangle = \{Q \in P(S|S \times A) | Q \leq \underline{Q} \leq \bar{Q}\}$ とするとき、 $\mathcal{Q}$ はLemma 2.1より凸多面体であるから、ある端点の集合 $\{Q^{(1)}, Q^{(2)}, \dots, Q^{(l)}\}$ によって $\mathcal{Q} = \text{conv}\{Q^{(1)}, Q^{(2)}, \dots, Q^{(l)}\}$ と表すことができる。 $\mathcal{Q} \ni Q = (q_{ij})$ について、各*i*行目ごとに推移確率行列の条件 $\sum_{j=1}^3 q_{ij} = 1$  ( $i = 1, 2, 3$ )をみたす端点調べる。 $\mathcal{Q}$ 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なようになる.

$$\text{ext}(\hat{q}_1) = \{(0.437, 0.187, 0.376), (0.4, 0.224, 0.376), (0.448, 0.26, 0.292), (0.489, 0.219, 0.292), (0.4, 0.26, 0.34), (0.489, 0.187, 0.324)\},$$

$$\text{ext}(\hat{q}_2) = \{(0.187, 0.437, 0.376), (0.224, 0.4, 0.376), (0.26, 0.448, 0.292), (0.219, 0.489, 0.292), (0.26, 0.4, 0.34), (0.187, 0.489, 0.324)\},$$

$$\text{ext}(\hat{q}_3) = \{(0.196, 0.262, 0.542), (0.168, 0.29, 0.542), (0.208, 0.334, 0.458), (0.235, 0.307, 0.458), (0.168, 0.334, 0.498), (0.235, 0.262, 0.503)\}$$
 を得る.

$\mathcal{Q}$  の端点の集合は,  $\hat{q}_1, \hat{q}_2, \hat{q}_3$  の端点からそれぞれ一つずつ選んで作る推移確率行列  $Q^{(l)}$  から成る ( $6^3$  個). 例えば,  $\text{ext}(\hat{q}_1)$  から成る凸多面体を図示すると次のような  $x+y+z=1$  の平面上の六角形となる (Figure 1).

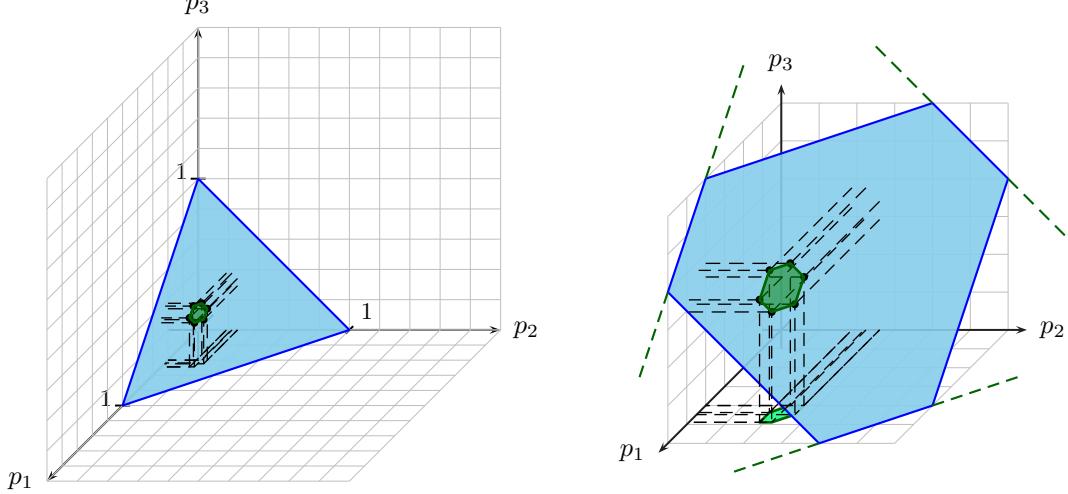


Figure 1: Graphs of interval  $\hat{q}_1$

$\beta = 0.9, \mathbf{r} = (3, 1, 2)', F \ni f$ (固定) として

$$\underline{L}(f)\mathbf{x} = \mathbf{r}(f) + \beta \min_{Q \in \mathcal{Q}(f)} Q\mathbf{x}$$

$$\overline{L}(f)\mathbf{x} = \mathbf{r}(f) + \beta \max_{Q \in \mathcal{Q}(f)} Q\mathbf{x}$$

の不動点を求めてみると,  $\underline{\phi}(f) = (20.003, 17.508, 18.643), \overline{\phi}(f) = (21.732, 19.232, 20.339)$  を得る. 従って, Teorem 3.1 から  $\phi(f|\mathcal{Q}(f)) = [\underline{\phi}(f), \overline{\phi}(f)]$  は次のように得られる:

$$\phi(f|\mathcal{Q}(f)) = ([20.003, 21.732], [17.508, 19.232], [18.643, 20, 339]).$$

真の推移確率行列を  $\begin{pmatrix} \frac{1}{2} & \frac{1}{6} & \frac{1}{3} \\ \frac{1}{3} & \frac{1}{3} & \frac{1}{3} \\ \frac{2}{5} & \frac{2}{5} & \frac{1}{5} \end{pmatrix}$  としたときの value function の値は  $\phi = (22.469, 20.116, 21.135)$  である.

## 8 区間ベイズ推定

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,  $\mathcal{Q} = \langle \underline{Q}, \overline{Q} \rangle \in \mathcal{M}_n$  の  $\underline{Q}, \overline{Q} \in \mathbb{R}_{+}^{n \times n}$  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 8.1.** (i)  $\underline{Q}_t \downarrow \underline{Q}, \overline{Q}_t \uparrow \overline{Q}$  ( $t \rightarrow \infty$ ),  $\langle \underline{Q}_t, \overline{Q}_t \rangle \neq \emptyset$  ( $t \geq 1$ ) とする. このとき,  $\langle \underline{Q}_t, \overline{Q}_t \rangle \xrightarrow{\rho} \langle \underline{Q}, \overline{Q} \rangle$  ( $t \rightarrow \infty$ )

(ii)  $\underline{Q}_t \uparrow \underline{Q}, \overline{Q}_t \downarrow \overline{Q}$  ( $t \rightarrow \infty$ ),  $\langle \underline{Q}, \overline{Q} \rangle \neq \emptyset$  ( $t \geq 1$ ) とする. このとき,  $\langle \underline{Q}_t, \overline{Q}_t \rangle \xrightarrow{\rho} \langle \underline{Q}, \overline{Q} \rangle$  ( $t \rightarrow \infty$ )

上の Lemma 8.1 を用いて次が示される.

**Theorem 8.1.**  $\underline{Q}_t \rightarrow \underline{Q}, \bar{Q}_t \rightarrow \bar{Q}$  ( $t \rightarrow \infty$ ),  $\mathcal{Q}_t := \langle \underline{Q}_t, \bar{Q}_t \rangle \neq \emptyset$  ( $t \geq 1$ ),  $\mathcal{Q} := \langle \underline{Q}, \bar{Q} \rangle$ .  
このとき, 次が成り立つ:

$$(i) \quad \mathcal{Q}_t \rightarrow \mathcal{Q} \quad (t \rightarrow \infty)$$

$$(ii) \quad \phi(f|\mathcal{Q}_t) \rightarrow \phi(f|\mathcal{Q}) \quad (t \rightarrow \infty) \quad (f \in F).$$

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

$$(8.1) \quad N_T(j|i, a, H_T) := \sum_{t=0}^{T-1} I_{\{X_t=i, \Delta_t=a, X_{t+1}=j\}} \quad (T \geq 1)$$

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  $\mathcal{Q}(H_T) = \langle Q(H_T), \bar{Q}(H_T) \rangle = [\underline{q}_{ij}(a|H_T), \bar{q}_{ij}(a|H_T)]$  とする. すなわち,

$$(8.2) \quad \underline{Q}(H_T) := (\underline{q}_{ij}(a|H_T) : i, j \in S, a \in A) \in \mathbb{R}_+^{n \times nk}$$

$$(8.3) \quad \bar{Q}(H_T) := (\bar{q}_{ij}(a|H_T) : i, j \in S, a \in A) \in \mathbb{R}_+^{n \times nk}$$

として,

$$(8.4) \quad \mathcal{Q}(H_T) = \langle \underline{Q}(H_T), \bar{Q}(H_T) \rangle$$

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 を区間ベイズ推定 MDPs{ $\mathcal{Q}(H_T)$ } と言う.

$$(8.5) \quad N_T(i, a|H_T) := \sum_{j \in S} N_T(j|i, a, H_T) \quad (i \in S, a \in A)$$

とおく.

区間ベイズの性質 ([1]) および Theorem 8.1 を用いて次の結果を得る.

**Theorem 8.2.**  $\{X_0, \Delta_0, X_1, \Delta_1, \dots\}$  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

$$N_T(i, a|H_T) \rightarrow \infty \quad (T \rightarrow \infty)$$

とする. このとき, 確率 1 で区間ベイズ推定 MDPs{ $\mathcal{Q}(H_T)$ } は MDPs{ $Q$ } に収束する, すなわち, 次が成り立つ.

$$(i) \quad \mathcal{Q}(H_T) \rightarrow \{Q\} \quad (T \rightarrow \infty)$$

$$(ii) \quad \phi(f|\mathcal{Q}(H_T)) \rightarrow \phi(f|Q) \quad (T \rightarrow \infty), \quad (f \in F).$$

## References

- [1] Lorraine De Robertis and J. A. Hartigan. Bayesian inference using intervals of measures. *Ann. Statist.*, 9(2):235–244, 1981.
- [2] Bharat Doshi and Steven E. Shreve. Strong consistency of a modified maximum likelihood estimator for controlled Markov chains. *J. Appl. Probab.*, 17(3):726–734, 1980.
- [3] Nagata Furukawa. Characterization of optimal policies in vector-valued Markovian decision processes. *Math. Oper. Res.*, 5(2):271–279, 1980.
- [4] Darald 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] 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.
- [7] Masami Kurano, Jinjie Song, Masanori Hosaka, and Youqiang Huang. Controlled Markov set-chains with discounting. *J. Appl. Probab.*, 35(2):293–302, 1998.
- [8] Masami Kurano, Masami Yasuda, and Jun-ichi Nakagami. Interval methods for uncertain Markov decision processes. In *Markov processes and controlled Markov chains (Changsha, 1999)*, pages 223–232. Kluwer Acad. Publ., Dordrecht, 2002.
- [9] K. Kuratowski. *Topology. Vol. I.* New edition, revised and augmented. Translated from the French by J. Jaworowski. Academic Press, New York, 1966.
- [10] P. Mandl. Estimation and control in Markov chains. *Advances in Appl. Probability*, 6:40–60, 1974.
- [11] J. J. Martin. *Bayesian decision problems and Markov chains*. Publications in Operations Research, No. 13. John Wiley & Sons Inc., New York, 1967.
- [12] Martin L. Puterman. *Markov decision processes: discrete stochastic dynamic programming*. John Wiley & Sons Inc., New York, 1994. A Wiley-Interscience Publication.
- [13] Eilon Solan. Continuity of the value of competitive Markov decision processes. *J. Theoret. Probab.*, 16(4):831–845 (2004), 2003.
- [14] K. M. van Hee. *Bayesian control of Markov chains*, volume 95 of *Mathematical Centre Tracts*. Mathematisch Centrum, Amsterdam, 1978.
- [15] Samuel S. Wilks. *Mathematical statistics*. A Wiley Publication in Mathematical Statistics. John Wiley & Sons Inc., New York, 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.