

区間型マルコフ決定モデルについて (On a Markov decision model with interval-valued transition matrices)

神奈川大学・理学部 堀口 正之 (Masayuki HORIGUCHI)
Faculty of Science, Kanagawa University

1 はじめに

本報告では、推移法則が未知の場合のマルコフ決定過程 (Markov decision Processes, MDP) について、推移法則を区間で推定し、推移確率行列の各成分が閉区間で表され、推移確率行列はその閉区間内に成分をもつ確率行列の集合として構成される区間型マルコフ連鎖とそのもとのマルコフ決定過程の定式化を行う。割引総期待利得の評価基準での価値関数の表現と、 α -パーセントイル型の価値関数の表現について考察を行う。

先行研究として、区間型のマルコフ集合連鎖モデルを考察した Kurano et. al (1999) や Hartfiel (1998) がある。これらの研究では、マルコフ連鎖を構成する推移確率行列の各成分が単一閉区間のそれぞれで表され、その条件下での推移確率行列の集合は、その区間内にある任意の確率行列からなる集合と捉えられる。そして、価値関数についても区間表現され、最適方程式とそのもとの最適性 (パレート最適性) によって最適解 (最適政策) が特徴づけられる。本報告では、これらについてその概略を述べ、推移法則未知のマルコフ決定モデルにおいて、その評価基準としてパーセントイル型について考察をする。

2 準備

本節では、割引平均期待利得を評価基準とするマルコフ決定過程についてまず述べて、その後、推移法則が各成分ごとに閉区間で表現されるマルコフ集合連鎖モデルでの最適化についてみていく。

マルコフ決定過程は、次のような4つの要素からなる $\{S, A, Q, r\}$:

$S = \{1, 2, \dots, n\}$ は n 個の有限要素をもつ状態空間であり、 $A = \{a_1, a_2, \dots, a_K\}$ は K 個の有限要素をもつ決定空間である。また、 $Q = (q_{ij}(a)) \in P(S|S \times A)$ は K 種の未知の推移確率行列からなるパラメータ空間とする。さらに、 $r = (r(i, a)) \in B_+(S \times A)$ は、状態 $i \in S$ で決定 $a \in A$ を選択したときに生じる利得を表す関数である。関数 $f: S \rightarrow A$ によって、各状態 $i \in S$ に対する決定 $f(i)$ を表す。 f は各期での確定的定常政策 (f, f, \dots) も表し、簡単に $f = (f, f, \dots)$ と表す。真の推移確率行列を Q とするとき、確定的定常政策 f の下での価値関数 $\phi: S \rightarrow \mathbf{R}$ は次式で与えられる。

$$(1) \quad \phi(f|Q) = \sum_{t=0}^{\infty} (\beta Q(f))^t r(f),$$

区間型の推移確率行列 $\langle \underline{A}, \bar{A} \rangle$ は、次式のような推移確率行列 (q_{ij}) の集合として定義される。

$$(2) \quad \langle \underline{A}, \bar{A} \rangle := \left\{ Q = (q_{ij}) \in \mathbf{R}_+^{m \times n} \mid \underline{a}_{ij} \leq q_{ij} \leq \bar{a}_{ij}, \right. \\ \left. q_{ij} \geq 0, \sum_{j=1}^n q_{ij} = 1 \ (1 \leq i \leq m, 1 \leq j \leq n) \right\}.$$

一般に, $\mathbf{R}^{m \times n}$ 上の半順序 \preceq と \prec を次のように定義する. 次のような 2 つの実数値行列 $\mathbf{R}^{m \times n} \ni A = (a_{ij}), B = (b_{ij})$ に対して,

$$(3) \quad \begin{cases} A \preceq B & \text{if } a_{ij} \leq b_{ij} \ (1 \leq i \leq m, 1 \leq j \leq n) \\ A \prec B & \text{if } A \preceq B \text{ and } A \neq B \end{cases}$$

とおく. そして, 任意の順序のついた行列 $A \preceq \bar{A}$ に対して, 区間型の推移確率行列 (\underline{A}, \bar{A}) が定義されるのである.

$n \times n$ 行列を成分に持つ区間型非負値行列の集合を \mathcal{M}_n と表す. すなわち,

$$(4) \quad \mathcal{M}_n = \{ \langle Q, \bar{Q} \rangle \mid \langle Q, \bar{Q} \rangle \neq \emptyset, Q \preceq \bar{Q}, Q, \bar{Q} \in \mathbf{R}_+^{n \times n} \}$$

ここで, $Q_1, Q_2 \in \mathcal{M}_n$ に対してその積を次のように定義する.

$$(5) \quad Q_1 Q_2 = \{ Q_1 Q_2 \mid Q_1 \in Q_1, Q_2 \in Q_2 \}$$

同一の区間型行列 $Q \in \mathcal{M}_n$ の k 回の積は次式によって帰納的に求められる.

$$(6) \quad Q^k = Q^{k-1} Q \ (k \geq 2).$$

非負実数値の集合 \mathbf{R}_+ におけるすべての有界かつ閉区間の集合を $C(\mathbf{R}_+)$ と表し, その n 次元列ベクトルをなす集合を $C(\mathbf{R}_+)^n$ と表す, すなわち,

$$(7) \quad C(\mathbf{R}_+)^n = \{ D = (D_1, D_2, \dots, D_n)' \mid D_i \in C(\mathbf{R}_+) \ (1 \leq i \leq n) \}$$

とする. ただし, ベクトル d の転置ベクトルの記号として d' を用いる.

区間型非負値行列における演算として, それぞれ次のように定義する. $D = (D_1, D_2, \dots, D_n)', E = (E_1, E_2, \dots, E_n)' \in C(\mathbf{R}_+)^n, h \in \mathbf{R}_+^n, \lambda \in \mathbf{R}_+$ に対して,

$$(8) \quad \begin{aligned} D + E &= \{ d + e \mid d \in D, e \in E \}, \\ h + D &= \{ h + d \mid d \in D \}, \\ \lambda D &= \{ \lambda d \mid d \in D \}. \end{aligned}$$

また, $D := [\underline{d}, \bar{d}] = ([\underline{d}_1, \bar{d}_1], [\underline{d}_2, \bar{d}_2], \dots, [\underline{d}_n, \bar{d}_n])' \in C(\mathbf{R}_+)^n$, where, $\underline{d} = (\underline{d}_1, \underline{d}_2, \dots, \underline{d}_n) \in \mathbf{R}_+^n, \bar{d} = (\bar{d}_1, \bar{d}_2, \dots, \bar{d}_n) \in \mathbf{R}_+^n$ とそれぞれおくと, 任意の $D = (D_1, D_2, \dots, D_n)' \in C(\mathbf{R}_+)^n$ と $G \subset \mathbf{R}_+^{1 \times n}$ に対してその積 GD を次のように定める.

$$(9) \quad GD = \{ gd \mid g = (g_1, g_2, \dots, g_n) \in G, d = (d_1, d_2, \dots, d_n)' \in D, \\ d_i \in D_i \ (1 \leq i \leq n) \}$$

この時, 次のような性質を得る.

Lemma 2.1 ([7, 16])

- (i) 任意の $Q \in \mathcal{M}_n$ は $\mathbf{R}^{n \times n}$ における凸多面体である.
- (ii) 任意のコンパクト部分集合 $G \subset \mathbf{R}_+^{1 \times n}$ と集合 $D \in C(\mathbf{R}_+)^n$ に対して, 積 $GD \in C(\mathbf{R}_+)$ が成り立つ.

さらにいくつかの半順序について定義する. $C(\mathbf{R}_+)$ 上の半順序 \preceq, \prec は次のように定められる. $[c_1, c_2], [d_1, d_2] \in C(\mathbf{R}_+)$ に対して,

$$\begin{cases} [c_1, c_2] \preceq [d_1, d_2] & \text{if } c_i \leq d_i \ (i = 1, 2), \\ [c_1, c_2] \prec [d_1, d_2] & \text{if } [c_1, c_2] \preceq [d_1, d_2] \text{ and } [c_1, c_2] \neq [d_1, d_2]. \end{cases}$$

また, $C(\mathbf{R}_+)^n$ 上の半順序 \preceq, \prec は上述の $C(\mathbf{R}_+)$ での半順序を用いて次のように定められる. 区間型のベクトル $\mathbf{v} = (v_1, v_2, \dots, v_n)'$ と $\mathbf{w} = (w_1, w_2, \dots, w_n)' \in C(\mathbf{R}_+)^n$ に対して,

$$\begin{cases} \mathbf{v} \preceq \mathbf{w} & \text{if } v_i \preceq w_i \ (1 \leq i \leq n), \\ \mathbf{v} \prec \mathbf{w} & \text{if } \mathbf{v} \preceq \mathbf{w} \text{ and } \mathbf{v} \neq \mathbf{w}. \end{cases}$$

3 区間型モデルでの評価関数

本節では, さらにいくつかの記法を準備して区間推定されたマルコフ決定モデルの定式化をする.

\mathbf{R}_+^n における優界かつ閉部分集合 D_1 と D_2 に対してハウスドルフ距離 ρ を次のように定義する.

$$(10) \quad \rho(D_1, D_2) := \max\left\{ \sup_{x \in D_1} \inf_{y \in D_2} \|x - y\|, \sup_{y \in D_2} \inf_{x \in D_1} \|x - y\| \right\}.$$

ただし, $\|\cdot\|$ は \mathbf{R}^n でのユークリッド距離である.

状態空間 $S = \{1, 2, \dots, n\}$ と決定空間 $A = \{1, 2, \dots, k\}$ に対して, 次のようにそれぞれの確率分布を定義する.

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

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

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

また, 有限要素の集合 D 上のすべての非負実数値関数の集合を $B_+(D)$ とおく. D ($n = \#D$) に対して, $B_+(D)$ は \mathbf{R}_+^n と同一視できる.

前述した下記の割引平均期待利得 ϕ に対して, つぎの Lemma が成り立つことがよく知られている (cf. Puterman(1994)).

$$(11) \quad \phi(f|Q) = \sum_{t=0}^{\infty} (\beta Q(f))^t \mathbf{r}(f)$$

確定的定常政策の集合を F とおき, $f \in F$ に対して, 写像 $L(f) : \mathbf{R}_+^n \rightarrow \mathbf{R}_+^n$ は次のように定義される.

$$(12) \quad L(f)\mathbf{x} = \mathbf{r}(f) + \beta Q(f)\mathbf{x}, \quad \mathbf{x} = (x_1, x_2, \dots, x_n)' \in \mathbf{R}_+^n.$$

このとき, つぎが成り立つ

Lemma 3.1 (cf. Puterman[22])

(i) $L(f)$ は単調増加かつ縮小写像である、すなわち、

$$\begin{aligned} \mathbf{x} \leq \mathbf{x}' \text{ ならば } L(f)\mathbf{x} \leq L(f)\mathbf{x}' \text{ が各成分ごとに成り立つ,} \\ \|L(f)\mathbf{x} - L(f)\mathbf{x}'\| \leq \beta\|\mathbf{x} - \mathbf{x}'\| \quad (\mathbf{x}, \mathbf{x}' \in \mathbf{R}_+^n), \end{aligned}$$

ただし、 $\|\cdot\|$ は sup-ノルムである。

(ii) $\phi(f|Q)$ は $L(f)$ ただ一つの不動点である、すなわち、任意の $\mathbf{x} \in \mathbf{R}_+^n$ に対して、

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

が成り立つ。

真の推移確率行列 Q によるマルコフ決定過程を $\{S, A, Q, \mathbf{r}\}$ と表すことにし、 Q を区間表現によって推定してその区間型推移確率行列を $\underline{Q} = \langle \underline{Q}, \overline{Q} \rangle$ とおく、ただし、この行列集合は

$$\begin{aligned} (13) \quad \underline{Q} &= (\underline{q}_{ij}(a) : i, j \in S, a \in A) \in \mathbf{R}_+^{kn \times n}, \\ \overline{Q} &= (\overline{q}_{ij}(a) : i, j \in S, a \in A) \in \mathbf{R}_+^{kn \times n}, \\ \mathcal{Q} &= \langle \underline{Q}, \overline{Q} \rangle = \{Q \in P(S|S \times A) \mid \underline{Q} \leq Q \leq \overline{Q}\}. \end{aligned}$$

によって表される。今、 \mathcal{Q} によって構成されるマルコフ決定過程を $\{S, A, \mathcal{Q}, \mathbf{r}\}$ と表し、区間によって表現されたマルコフ過程と呼ぶことにしよう。このとき、評価関数も区間によって表現され次のように表される。

$$(14) \quad \phi(f|\mathcal{Q}) = \{\phi(f|Q) \mid Q \in \mathcal{Q}\} \subset \mathbf{R}_+^n$$

ただし、 $\phi(f|Q)$ は通常の推移確率行列からなるマルコフ決定過程での価値関数(式(11))である。次に、 $\phi(f|\mathcal{Q})$ の性質についてまとめておく。

まず、 $\phi(f|\mathcal{Q}) \in C(\mathbf{R}_+)^n$ が成り立つ。写像 $\mathcal{L} : C(\mathbf{R}_+)^n \rightarrow C(\mathbf{R}_+)^n$ を次のように定義する。

$$(15) \quad \mathcal{L}(f)\mathbf{v} = \mathbf{r}(f) + \beta \mathcal{Q}(f)\mathbf{v}, \quad \mathbf{v} \in C(\mathbf{R}_+)^n,$$

ただし $\mathcal{Q}(f) = \langle \underline{Q}(f), \overline{Q}(f) \rangle$, $\underline{Q}(f) = (\underline{q}_{ij}(f(i))) \in \mathbf{R}_+^{n \times n}$, $\overline{Q}(f) = (\overline{q}_{ij}(f(i))) \in \mathbf{R}_+^{n \times n}$. このとき Lemma 2.1 から、 $\mathcal{L}(f)\mathbf{v} \in C(\mathbf{R}_+)^n$ ($\mathbf{v} \in C(\mathbf{R}_+)^n$) であることもわかる。また、 $\underline{L}(f) : \mathbf{R}_+^n \rightarrow \mathbf{R}_+^n$, $\overline{L}(f) : \mathbf{R}_+^n \rightarrow \mathbf{R}_+^n$ とそれぞれ定義すると、 $\mathbf{x} = (x_1, x_2, \dots, x_n)' \in \mathbf{R}_+^n$ に対して以下の等式がそれぞれ成り立つ。

$$(16) \quad \underline{L}(f)\mathbf{x} = \mathbf{r}(f) + \beta \min_{Q \in \mathcal{Q}(f)} Q\mathbf{x}$$

$$(17) \quad \overline{L}(f)\mathbf{x} = \mathbf{r}(f) + \beta \max_{Q \in \mathcal{Q}(f)} Q\mathbf{x}$$

よって、以下を得る。

Lemma 3.2 任意の確定的定常政策 $f \in F$ に対して、

(i) $\mathcal{L}(f)$ は単調増加かつ縮小写像である。

(ii) $\underline{L}(f)$ と $\overline{L}(f)$ はともに sup-ノルムに関して単調増加かつ縮小写像である。

Lemma 3.1 と Lemma 3.2 から次を得る。

Theorem 3.1 任意の確定的定常政策 $f \in F$ に対して,

(i) $\phi(f|Q) \in C(\mathbf{R}_+)^n$ であって $\phi(f|Q)$ は $\mathcal{L}(f)$ のただ一つの不動点である. さらに, 任意の $v \in C(\mathbf{R}_+)^n$, に対して

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

が成り立つ.

(ii) $\phi(f|Q) = [\underline{\phi}(f), \overline{\phi}(f)]$ とおく. このとき, $\underline{\phi}(f)$ と $\overline{\phi}(f)$ はそれぞれ $\underline{\mathcal{L}}(f)$ と $\overline{\mathcal{L}}(f)$ の不動点である.

$u \in C(\mathbf{R}_+)^n$ に対して,

$$(18) \quad \mathcal{L}(u) := (\mathcal{L}(u)_1, \mathcal{L}(u)_2, \dots, \mathcal{L}(u)_n)',$$

とおく, ただし $\mathcal{L}(u)_i := \text{eff}(\{r(i, a) + \beta Q_{i,a} u | a \in A\})$ ($i \in S$) であってこれは有効点の集合を表す.

Lemma 3.3 2つの確定的定常政策 $f, g \in F$ に対して, もし $\phi(f|Q) \prec \mathcal{L}(g)\phi(f|Q)$ が成り立つとき,

$$\phi(f|Q) \prec \phi(g|Q)$$

を得る.

Theorem 3.2 $f^* \in F$ がパレート最適であるための必要かつ十分条件は次の最適包含関係の最大解であることである

$$(19) \quad u \in \mathcal{L}(u), u \in C(\mathbf{R}_+)^n.$$

4 区間推定されたMDPの考察

ここでは, 区間推定されたマルコフ決定過程の構成とパーセントイル型の評価について考察する.

De Robertis と Hartigan のベイズの区間推定手法の結果 ([25]) を, 我々の推移確率行列未知のマルコフ決定過程での真の推移確率行列の区間推定表現に適用する. 具体的には, 推移確率行列の各行において, その行ベクトルの成分の推定を行う. よって, $P_n := P(S) = \{p = (p_1, p_2, \dots, p_n) | p_i \geq 0, \sum_{i=1}^n p_i = 1\}$ とおき, \mathcal{B} を \mathbf{R}^N のすべての可測集合の全体とおく. 通常のベイズ推定では, 事前分布として事後分布の表現に便利な分布族を用いるが, ここでは, 事前分布として, \mathcal{B} 上の区間で表現可能な測度の集合を事前分布として扱う. それらは, 次のようにして得られる. \mathcal{B} 上の2つの測度 L と U がすべての可測集合 $A \in \mathcal{B}$ に対して $L(A) \leq U(A)$ であるとき, 単に $L \leq U$ と表すことにする. また, $[L, U]$ はすべての $A \in \mathcal{B}$ に対して $L(A) \leq Q(A) \leq U(A)$ を満たす測度 Q の凸集合であって, これを事前測度として扱う.

以下では, 議論の単純化のため $[L, kL]$ ($k \geq 1$) を仮定する, ただし $L(\cdot)$ は P_n 上のルベーク測度である.

事前測度 $[L, U]$ によって, 推移確率行列の各成分 p_i は次のような区間 $[\underline{\lambda}_i, \overline{\lambda}_i]$ ($i \in S$) によって事後推定され, とくにこの両端点は次の積分比の区間表現として得られる.

$$(20) \quad \left\{ \int_{P_n} p_i Q(dp) / \int_{P_n} Q(dp) \middle| L_\sigma \leq Q \leq U_\sigma \right\},$$

ただし, L_σ と U_σ はそれぞれ観測データ $\sigma = (\sigma_1, \sigma_2, \dots, \sigma_n)$ に対する事後測度の下限と上限であり, 以下の方程式 (21) と (22) のただ一つの解にそれぞれなっている.

Theorem 4.1 下限 $\underline{\lambda}_i$ と上限 $\bar{\lambda}_i$ はそれぞれ次の方程式のただ一つの解である.

$$(21) \quad U_\sigma(p_i - \underline{\lambda}_i)^- + L_\sigma(p_i - \underline{\lambda}_i)^+ = 0,$$

$$(22) \quad U_\sigma(p_i - \bar{\lambda}_i)^+ + L_\sigma(p_i - \bar{\lambda}_i)^- = 0,$$

ただし, $Q(f)$ は可測関数 f の測度 Q に関する積分値を表し, $x^+ = \max\{0, x\}$, $x^- = x - x^+ = \min\{0, x\}$ である.

α -パーセンタイル型の区間推定に関しては次のようになる. $\underline{g}_{i,a}$ と $\bar{g}_{i,a}$ を P_n 上の次のような可測関数とおく.

$$(23) \quad \underline{g}_{i,a}(p) = I_{\{p_i \leq a\}}(p), \quad \bar{g}_{i,a}(p) = I_{\{p_i \geq a\}}(p),$$

ただし, I_A は A の指示関数で $I_A(x) = 1$ if $x \in A$, $= 0$ if $x \notin A$ である. また,

$$(24) \quad \underline{\lambda}(a|\sigma) = \sup \left\{ \frac{Q_\sigma(\underline{g}_{i,a})}{Q_\sigma(I_{P_n})} \mid Q_\sigma \in [L_\sigma, kL_\sigma] \right\},$$

$$(25) \quad \bar{\lambda}(a|\sigma) = \sup \left\{ \frac{Q_\sigma(\bar{g}_{i,a})}{Q_\sigma(I_{P_n})} \mid Q_\sigma \in [L_\sigma, kL_\sigma] \right\}.$$

とおく. $\underline{p}_i(\alpha)$ と $\bar{p}_i(\alpha)$ を次のように定義する.

$$(26) \quad \underline{\lambda}(\underline{p}_i(\alpha)|\sigma) = \alpha, \quad \bar{\lambda}(\bar{p}_i(\alpha)|\sigma) = \alpha.$$

このとき, $\underline{p}_i(\alpha)$ と $\bar{p}_i(\alpha)$ は次を満たす. $B(\alpha, \beta)$ はベータ関数, $B(\alpha, \beta|x)$ は不完全ベータ関数を表す.

Theorem 4.2 $\underline{p}_i(\alpha)$ と $\bar{p}_i(\alpha)$ について

$$(27) \quad \frac{B(s, t|\underline{p}_i(\alpha))}{B(s, t)} = \frac{\alpha}{\alpha + (1 - \alpha)k}, \quad \frac{B(s, t|\bar{p}_i(\alpha))}{B(s, t)} = \frac{(1 - \alpha)k}{\alpha + (1 - \alpha)k}$$

が成り立つ.

例えば, 各行ごとの状態推移の回数を数えて次のような観測結果 Σ を得たとき

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

このデータセットから, $\alpha = 0.05$ としたときの区間推定された推移確率行列は

$$\begin{pmatrix} [0.095, 0.837] & [0.013, 0.641] & [0.045, 0.750] \\ [0.013, 0.641] & [0.095, 0.837] & [0.045, 0.750] \\ [0.011, 0.595] & [0.039, 0.701] & [0.140, 0.860] \end{pmatrix}$$

を得る. また, これをもとにした価値関数の区間表現は

$$\phi(f|Q(f)) = ([14.0878, 27.9643], [11.8025, 25.66], [12.9396, 26.656])$$

である。

[α -パーセンタイル型評価の課題]

真の推移確率行列が既知の場合には、value $\phi(f|\mathcal{Q})$ の閾値評価問題として次のような考察が可能である (cf. J.Filar, D.Krass, K.W. Ross (1995) など) :

$$\text{Prob}_u(\phi \geq \tau | X_1 = s_1) \geq \alpha$$

を満たす価値関数の値が τ 以上であることの政策 u の存在性の議論を行い、最適解ではなくあるレベル α 以上の政策 u の存在において最適化問題を議論する。これらでは、Occupation measure や parametric LP での考察が最適政策を導出可能かさらに検討が必要となる。一方、推移法則未知のモデルでの α -percentile の意味は、統計数学での検出力に関するアプローチでもあり、価値関数の区間表現とそれをもとにした意思決定のもつ意味についてさらに検討が必要で、これらが今後の課題の一つである。

References

- [1] Bertsekas, D.P. and Shreve, S.E. (1978). *Stochastic Optimal Control*. New York: Academic Press.
- [2] Dynkin, E.B. and Yushkevich A.A. (1979). *Controlled Markov Processes and their Applications*. New York - Berlin: Springer-Verlag.
- [3] Ferguson, T.S. (1967). *Mathematical Statistics*. New York - London: Academic Press.
- [4] DeGroot, M.H. (1970). *Optimal Statistical Decisions*. New York: McGraw-Hill Book Co.
- [5] Furukawa, N. and Iwamoto, S. (1970). Stopped decision processes on complete separable metric spaces. *J. Math. Anal. Appl.*, 31:615–658.
- [6] van Hee, K.M. (1978). *Bayesian Control of Markov Chains*. Mathematical Centre Tracts, No. 95. Amsterdam: Mathematisch Centrum.
- [7] Darald J. Hartfiel. *Markov set-chains*, volume 1695 of *Lecture Notes in Mathematics*. Springer-Verlag, Berlin, 1998.
- [8] Hernández-Lerma, O. and Marcus, S.I. (1985). Adaptive control of discounted Markov decision chains. *Journal of Optimization Theory and Applications* 46: 227–235.
- [9] Hernández-Lerma, O. (1989). *Adaptive Markov Control Processes*, volume 79 of *Applied Mathematical Sciences*. New York: Springer-Verlag.
- [10] Hernandez-Lerma, O. and Lasserre, J.B. (1996). *Discrete-time Markov Control Processes*. New York: Springer.
- [11] Hordijk, A. (1974). *Dynamic Programming and Markov Potential Theory*. Mathematical Centre Tracts, No. 51. Amsterdam: Mathematisch Centrum.
- [12] Horiguchi, M. (2015). Adaptive methods for multivariate Bayesian Control Chart (II), *RIMS kokyuroku No. 1939 (In Japanese)*, 152–161.

- [13] Horiguchi, M. (2015). Bayesian Inference in Markov Decision Processes. In, *Modern Trends in Controlled Stochastic Processes: Theory and Applications, Vol. 2 (A.B. Piunovskiy ed.)*, Luniver Press, 177–189.
- [14] Horiguchi, M. and Piunovskiy, A.B. (2013). Optimal stopping model with unknown transition probabilities. *Control Cybernet.*, 42(3):593–612.
- [15] Iki, T., Horiguchi, M., Yasuda, M. and Kurano, M. (2009). An interval Bayesian Method for uncertain MDPs.(In Japanese), *RIMS Kokyuroku* 1636:1–8.
- [16] Masami Kurano, Jinjie Song, Masanori Hosaka, and Youqiang Huang. Controlled Markov set-chains with discounting. *J. Appl. Probab.*, 35(2):293–302, 1998.
- [17] Kurano, M. (1972). Discrete-time Markovian decision processes with an unknown parameter. Average return criterion. *Journal of the Operations Research Society of Japan* 15: 67–76.
- [18] Kurano, M. (1983). Adaptive policies in Markov decision processes with uncertain transition matrices. *Journal of Information and Optimization Sciences* 4: 21–40.
- [19] Makis, V. (2008). Multivariate Bayesian control chart. *Oper. Res.*, 56(2):487–496, 2008.
- [20] Mandl, P. (1974). Estimation and control in Markov chains. *Advances in Applied Probability* 6: 40–60.
- [21] Martin, J.J. (1967). *Bayesian Decision Problems and Markov Chains*. Publications in Operations Research, No. 13. New York: John Wiley & Sons Inc.
- [22] Martin L. Puterman. *Markov decision processes: discrete stochastic dynamic programming*. John Wiley & Sons Inc., New York, 1994. A Wiley-Interscience Publication.
- [23] Raiffa, H and Schlaifer, R. (1961). *Applied Statistical Decision Theory*. Studies in Managerial Economics. Division of Research, Graduate School of Business Administration, Harvard University, Boston, Mass.
- [24] Rieder, U. (1975). Bayesian Dynamic Programming. *Advances in Applied Probability* 7: 330–348.
- [25] Lorraine De Robertis and J. A. Hartigan. Bayesian inference using intervals of measures. *Ann. Statist.*, 9(2):235–244, 1981.
- [26] Ross, S.M. (1970). *Applied Probability Models with Optimization Applications*. San Francisco: Holden-Day.
- [27] Ross, S.M. (1983). *Introduction to Stochastic Dynamic Programming*. San Diego, CA: Academic Press.
- [28] Wald, A. (1950). *Statistical Decision Functions*. New York: John Wiley & Sons Inc.
- [29] White, D.J. (1969). *Dynamic Programming*. Mathematical Economic Texts, 1. Edinburgh-London: Oliver & Boyd.