

適応制御過程

歳野 正美

1. はじめに

個人あるいは団体の意思決定は、ほとんどの場合未来に向けての行動選択と考えられる。未来における結果は本質的に確定的でありえないから、不確実性のもとでの決定問題を考えねばならない。たとえば幼稚園児の遠足では、ご馳走弁当の下ごしらえ等を前の晩に用意する。しかし、もし当日雨が降って遠足が中止になるとご馳走がふいになってしまう。これは意思決定の損失と考えられよう。逆の場合もありえる。すなわち、明日は雨だと思って、弁当の用意をしなかったが、当日雨がやんで遠足が可能になった場合、園児は落胆し主婦の嘆きはさぞ大変なものだろう。賢い主婦は明日雨が降るか否かを予測してご馳走を作るかどうか、損失を考えてその内容および経費の掛け具合を決めるであろう。このように、われわれの行動決定には多くの場合不確実性のもとでなされるが、これに対処する1つの方法として、不確実性を確率（法則）で測り、各行動に対する効用の期待値を最大にするように行動選択を行なうことが考えられる[13]。しかし、この場合不確実性のありようの表現である確率法則をいかにして定めるかが大問題である。遠足の例では、テレビの天気予報を見るとか、友人の意見を聞くとかしてよりよい決定を行なうためにいろいろと明日の天気についての情報を集めるであろう。新しい情報が

はいるたびに、さらに新しい情報を得るために行動をおこすべきか、もしおこすとしたらどのような行動をとるべきか、あるいは、手持ちの情報で満足して最終的な判断を下すべきかと主婦の気持は動揺して、いくつかの屈折を経たのち最終決定がなされる。この身近な例でもわかるように、未知の要素が含まれている最適化問題では未知の要素についての有益な情報をいかにして収集し、また得られた情報を決定過程の中に組み込み、いかにして最適化をはかってゆくかの問題がとり扱われる。いわゆる**適応型**の問題といわれるものである。このように考えてくると、ORでとり扱われる数学モデル、たとえば在庫モデル、取替モデル、ゲームモデルなどにはすべて**適応型**のモデルのあてはまる場合が非常に多い。しかもこれは現実をより反映したモデルといえよう。ところが**適応型**の問題は情報をとり扱うため、一般的にモデルが複雑になり解析的にも数値的にも最適解を求めるのが困難である。これに対処する一方法として、ベルマン (Bellman) の開発した**DP** [1] は最有力である。DPの手法があるからこそ**適応型**の問題が解けるようになったと言え過ぎるであろうか。ここでは、**適応型**の問題を**適応制御過程**あるいは**マルコフ決定過程**として定式化し、DPの方法を使って解く方法を解説する。さらに応用範囲の広い**マルコフ過程**の**適応制御**をとりあげて情報処理と最適化とに関するいくつかのアプローチの方法と結果を紹介する。

2. マルコフ決定過程と DP

たばこ製造会社は年度末の自社製品のシェアにもとづいて次の年度の広告費を決める。ある計画期間(会計年度を1期間として、たとえば N 期間)にわたって総利益を最大にするためには、各年度にどのくらいの広告費を支出すればよいであろうか。タバコのシェアは年々変動するがその年度末のシェアを変化する状態と考え、支出される広告費を行動 (action) としてこの問題は逐次決定過程、特にマルコフ決定過程 (Markov Decision Process, **MDP**) によって定式化される。その前にMDPの定義を与えておく。 n 期でのシステムの状態 i_n は状態空間 $S = \{1, 2, \dots, L\}$ の1点として表わされ、行動 a_n は行動空間 $A = \{1, 2, \dots, K\}$ から決定者によって選択されるものとする。たとえば、たばこ製造会社の問題ではシェアが15%以上(未満)である状態を1(2)で表わせれば、 $S = \{1, 2\}$ となる。状態の推移 (dynamics) に関しては定常性とマルコフ性を仮定する。すなわち、推移確率行列 $q = (q_{ij}^k)$ 、ただし $q_{ij}^k \geq 0$ 、 $\sum_{j=1}^L q_{ij}^k = 1$ が与えられていて、 $\text{Prob}(i_{n+1}=j | i_0, a_0, \dots, i_n=i, a_n=k) = q_{ij}^k$ 、 $i, j \in S, k \in A$ が成り立つものとする。また与えられた利得行列 $r = (r(i, k))$ に対して、 $i_n=i$ で行動 $a_n=k$ を選ぶと利得 $r(i, k)$ が発生する。各期の行動 a_n を選択する規則を政策という。一般に a_n は n 期以前の履歴 $h_n = (i_0, a_0, \dots, i_n)$ に依存して確率的に選ばれるが、特に n 期の状態 i_n のみに依存して確率1である行動を選ぶ政策を確定的なマルコフ型の政策といい、 S から A への写像 f_n の組 $\pi = (f_0, f_1, \dots, f_N)$ で表わされる。この政策を使うと n 期では $i_n=i$ のとき確率1で $i_n=f_n(i)$ が選択されることになる。また $f = f_n(n \geq 0)$ のとき定常政策といい、簡単のために f で表わす。政策 π を使ったときの総期待利得は初期状態 $i_0=i$ の関数として次のように表わされる

$$(2.1) \quad V_N(i, \pi) = \sum_{n=0}^N E_x[\beta^n r(i_n, a_n) | i_0=i].$$

ただし、 $0 < \beta \leq 1$ は割引き率。

(2.1)を最大にする最適政策を求める問題を考

えてみよう。目的関数(2.1)には加法性があるのでDPの最適性の原理[1]を適用して解くことができる。

$V_n(i) = n$ 期間問題で初期状態が i のとき、最適政策のもとでの最大総期待利得とすると次の再帰関係式を満足する。

$$(2.2) \quad V_0(i) = \max_{k \in A} r(i, k)$$

$$V_n(i) = \max_{k \in A} \{r(i, k) + \beta \sum_{j=1}^L q_{ij}^k V_{n-1}(j)\}$$

$$i \in S, n=1, 2, \dots, N$$

ここで、各 $i \in S$ に対して(2.2)の右辺の最大値を与える1つの行動を $f_n(i)$ として、 $f_n: S \rightarrow A$ を定義すると $\pi^* = (f_N, f_{N-1}, \dots, f_0)$ は最適政策となる。

無限期間($N = \infty$)での基準として

(i)割引き率 β ($0 < \beta < 1$)のある割合(割引き最適基準)

$$(2.3) \quad V_\beta(i, \pi) = \sum_{n=0}^{\infty} E_x[\beta^n r(i_n, a_n) | i_0=i]$$

(ii)長期間における1期当りの平均期待利得(平均最適基準)

$$(2.4) \quad g(i, \pi) = \lim_{N \rightarrow \infty} \inf (N+1)^{-1} \sum_{n=0}^N E_x[r(i_n, a_n) | i_0=i]$$

たとえば、表1で示すデータのもとでの平均最適基準における最適政策は状態1では2 (low advertising)、状態2では1 (high advertising)をとる定常政策となり、そのときの1期当りの平均利得=56.37となる。

MDPはその構造の一般性により、ORでとり扱う決定過程、統計学の逐次解析、最適制御など応用範囲がきわめて広く、この20年間で膨大な量の研究論文が発表されている。MDPの一般理論、計算アルゴリズム、応用に関する survey はそれぞれ[4, 14, 9]に与えられている。またMDPの日本語の教科書として[5, 6, 18]をあげておく。

3. 適応型の決定モデル

今、2台のスロットマシンI、IIがある。マシンIを使えば確率 r で1ドルが得られ、マシンII

表 1 たばこ製造会社のデータ

状態 i	行動 k	推移確率		利得
		q_{1k}	q_{2k}	r_k
1. シェアが 15%以上	1	0.7	0.3	68
	2	0.5	0.5	80
2. シェアが 15%未満	1	0.6	0.4	25
	2	0.33	0.67	33.2

ただし、行動 1 (2) : 高い(安い)広告費の支出

を使えば確率 p で 1 ドルが得られるとする。各期ではマシン I, II のどちらかを選んで使う。このとき、 N 期間に得られる総期待金額を最大にするためには、各期でマシン I, II のどちらを選んで使えば良いであろうか。もし r と p の値が既知であれば、1 ドルが出てくる確率 r 、 p の大きいほうのマシンを每期使えばよい。しかし、 r 、 p の値が未知の場合はどうふるまえばよいであろうか。これは逐次実験計画の中で特に **2 腕の盗賊の問題 (two-armed bandit problem)** [2, 17, 20] と呼ばれている。今、簡単のために r の値は既知で p の値が未知であるとしよう。 p の値が未知であっても、第 1 期目の選択を行なう時に決定者は p に対するなんらかの情報・予備知識 (これを**初期情報**という) をもっていることが考えられる。たとえば、マシン II を使ったことのある友人から情報の提供を受けたり、宣伝文を読んだりして初期情報を豊富にするように努めるかもしれない。またなんの予備知識をもち合せてない場合でも広く解釈して、それ自体がまた 1 つの初期情報と考えられる。今、マシン II を最初から n 期のあいだ、つづけて使ったとする。このとき、 n 回の実験結果 $1100 \cdots 10$ から得られる情報が初期情報に追加され、これらの情報をもとにして $n+1$ 期の選択がなされる。ただし 1 (0) は成功(失敗)を表わす。

一般に未知パラメータを含む決定過程を DP で定式化して解くためには、決定に役立つあらゆる場合の情報がある集合 \mathcal{S}_1 の一点として表わされかつ過程の進行にともない追加される情報が \mathcal{S}_1 上の変換あるいは推移として表わされる必要が

ある。情報をこのような形で表わしたものを**情報様式 (Information Pattern)** という。この場合 \mathcal{S}_1 は過程の状態空間の一部を構成し、追加情報による情報様式の改訂の方法が状態の推移(確率)法則を規定することになる。

以上を考慮して典型的な**適応制御過程**を定式化してみる。過程の状態はその期までの系の状態 p とその期までの未知パラメータに関する情報の表現である情報様式 P の対 (p, P) で表わされ、 (p, P) のとりうる値の全体を \mathcal{S} とする。状態 $(p, P) \in \mathcal{S}$ に対して行動空間 \mathcal{A} から 1 つの行動 $a \in \mathcal{A}$ を選択すると利得 $R((p, P), a)$ が発生する。またこのとき系の状態の推移と追加情報が得られるがこのメカニズムが状態空間 \mathcal{S} の上の推移確率法則 $Q(\cdot | (p, P), a)$ で表わされると仮定する。

明らかに上で定義された決定過程は \mathcal{S} を状態空間にもつ MDP であり、DP の再帰関係式 (2.2) を解くことによって最適な適応政策が求められる。現実の具体的問題に対しては、情報様式を選択およびその改訂の方法が問題になるが、これには統計学の推測論、特に十分統計量が利用される。通常、未知パラメータの“もっともらしさ”の程度を表わす確率分布で情報を記述して、追加情報をベイズの定理によって事前分布から事後分布に改定する方法がよく用いられる。これを**ベイズ (Bayes) モデル**と呼ぶ。スロットマシンの問題を Bayes モデルで解いたのが図 1 である。ただし $N=6$ 、 $r=0.6$ 、初期情報は一様分布。このとき最大期待利得は 3.72 ドルとなり、常に I を使うときの期待利得 3.6 ドルより 0.12 ドルと高くなっている。2 腕の盗賊の問題については、医学的な治療法の選択の問題などに応用されている ([15])。

4. マルコフ過程の適応制御

適応型の逐次決定モデルでは、情報様式を選択およびその改訂の方法が重要であることはすでに述べたが、そのアプローチの方法の違いによって、いろいろなタイプの適応制御モデルが考えられ

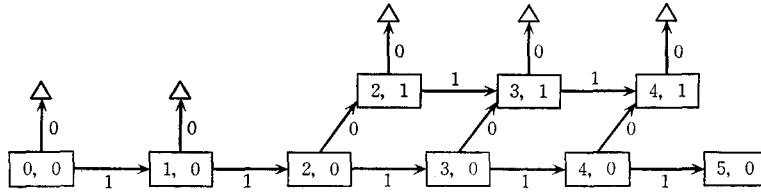


図 1 最適政策 $N=6, r=0.6$

ただし $[i, j]$ は Π を使用, $i(j)$ は成功(失敗)の回数, Δ では以後 I のみ使用

る. ここでは, MDP で特に状態の推移確率法則に未知のパラメータが含まれる場合, すなわち, m 次元ユークリッド空間の部分集合 θ をパラメータ空間として, 各 $\theta \in \Theta$ に対して 1つの MDP の推移確率法則 $q(\theta) = (q_{ij}^k(\theta))$ が対応しており, θ の値が未知の場合についての適応制御について考えてみよう. この場合, (2.3)あるいは(2.4)で定義された V_β, g は未知パラメータの値 θ に依存するので, θ の関数として, $V_\beta(i, \pi, \theta), g(i, \pi, \theta)$ と表わされる.

4.1 ベイズモデル

すでに述べたように, 初期情報をパラメータ空間 θ の上の事前確率分布で表わし, 情報の蓄積を事後分布で記述する方法がベイズモデルである. 相当する適応制御過程の n 期の状態は i_n と n 期の履歴のもとでの θ の上の事後分布 ξ_n の対 (i_n, ξ_n) で表わされ, したがって状態空間は $S \times P(\theta)$ となる. ただし $P(\theta)$ は θ の上の確率分布の全体. 事前分布 $\xi \in P(\theta)$ に対して, 状態 i で行動 k を選択し, 次の期に状態が j に移行したという情報のもとでの事後分布を $T_{ij}^k \xi$ と表わすと, DP の再帰関係式(2.2)は次のようになる. $n \geq 1, (i, \xi) \in S \times P(\theta)$ に対して

$$(4.1) \quad V_n(i, \xi) = \max_{k \in A} U_k V_{n-1}(i, \xi),$$

$$\text{ただし } V_0(i, \xi) = \max_{k \in A} r(i, k)$$

$$U_k u(i, \xi) = r(i, k) + \beta \sum_{j \in S} q_{ij}^k(\theta) u(j, T_{ij}^k \xi) \xi(d\theta)$$

(4.1)で $V_\beta(i, \xi) = \lim_{n \rightarrow \infty} V_n(i, \xi)$ (ただし $0 < \beta < 1$) とおけば, 割引最適基準の最適方程式

$$(4.2) \quad V_\beta(i, \xi) = \max_{k \in A} U_k V_\beta(i, \xi) \text{ を得る.}$$

今, $f(i, \xi) = \arg \max_{k \in A} U_k V_\beta(i, \xi)$ とする時, n

期の行動として $a_n = f(i_n, \xi_n)$ を選ぶ定常政策が **ベイズ最適** となることが知られている. すなわち

$$V_\beta(i, \xi) = \max_{\pi} \int V_\beta(i, \pi, \theta) \xi(d\theta) \\ = \int V_\beta(i, f, \theta) \xi(d\theta).$$

しかし, ベイズ最適政策を求めめるためには, 結局(4.2)式を解く必要があり困難が予想される. そこで, 簡単でとり扱いやすい Bayesian equivalent rule などの近似的な政策が提案されその性質が調べられている[19].

次に未知パラメータに関する情報の収集という立場からベイズ最適政策を検討してみよう.

例 ([8]) $S=A=\{1, 2\}, r_1^1=r_1^2=1, r_2^1=r_2^2=0$, 推移行列 q は図2のように, $q_{22}^2 = \theta = 1 - q_{21}^2$ のみが未知で他は既知とする. この場合 $\theta = [0, 1]$ で, 初期分布を適当にとるとベイズ最適政策のもとでは, 確率 1 で $a_n = 1 (n \geq 0)$ となることが最適方程式(4.2)から示される. これは状態 2 において行動 1 を常にとることを意味しており, 未知パラメータ θ に対してなんらの情報も得られないことになる. そのために, たとえば θ が 0 に十分近いとき明らかに状態 2 では行動 2 をとるのが最適であるのにベイズ政策のもとでは 2 を選ばないで 1 を選ぶという結果になっている.

このようなベイズ政策の欠点は漸近的最適性の概念, あるいは Forced Choice Circle などの手法が導入されて解決されている ([8, 16, 19]).

4.2 ノンベイズモデル

事前分布の存在を仮定しない, いわゆる Non-Bayesian 的な方法を考察してみよう. 各 $\theta \in \Theta$ および任意の政策 π と初期状態 $i \in S$ に対して

$$(4.3) \quad g(i, \pi^*, \theta) \geq g(i, \pi, \theta)$$

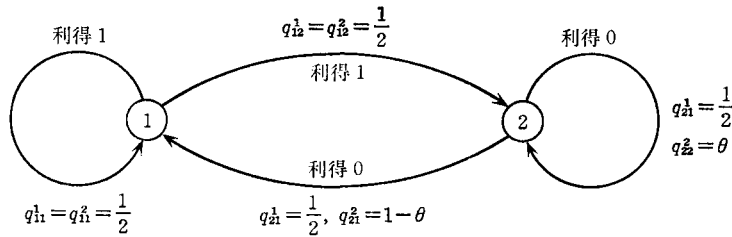


図 2 未知パラメータをもつ推移確率行列

が成り立つ π^* を平均最適な適応政策と呼び、 π^* の構成法の問題を検討してみよう。1つの考え方として、蓄積された情報をもとにして未知のパラメータを推定し、この推定値がパラメータの真値とみなしたとき、最適となる行動を選択する方法、すなわち**推定と制御の原理** ([7, 12]) を挙げることができる。今、この考え方によって1つの適応政策を構成してみよう。推定量として最尤推定量 (MLE) を利用することにする。

$$(4.4) \quad \hat{\theta}_n \equiv \arg \max_{\theta \in \Theta} L_n(\theta, h_n) \quad (n \geq 0)$$

$$\text{ここで } L_n(\theta, h_n) = \sum_{i=0}^{n-1} \log q(i_{i+1} | i_i, a_i)$$

$$\text{ただし } q(j | i, k) = q_{ij}^k, \quad h_n = (i_0, a_0, \dots, i_n)$$

今、各 $\theta \in \Theta$ に対して $g(i, f[\theta], \theta) = \sup_{\pi} g(i, \pi, \theta)$ が成り立つ最適定常政策 $f[\theta]: S \rightarrow A$ が存在すると仮定しよう。このとき n 期において履歴 h_n から (4.4) によって $\hat{\theta}_n$ を求め、 $\hat{\theta}_n$ をパラメータの値とみなしたときに最適となる行動 $a_n = f[\hat{\theta}_n](i)$ を選択する政策 $\pi^* = (f[\hat{\theta}_1], [\hat{\theta}_2], \dots)$ が考えられる。この適応政策 π^* を使った時ある条件のもとで MLE の一致性、すなわち確率 1 で $\hat{\theta}_n$ がパラメータの値 θ_0 に収束すること、および π^* が平均最適となることが示されている ([7, 12])。

4.3 学習モデル

各期に得られる情報をいわゆる reward-penalty 型の学習方式 ([11]) と value-iteration ([3]) で処理する方法について述べよう。 n 期での行動 a_n に対する条件付き確率分布を

$$\pi_n(k | i) = \text{Prob}(a_n = k | i_0, a_0, \dots, i_n = i) \quad (k \in A)$$

で表わす。このとき π_n が n を増大させるといくらかでも平均最適な政策に近づく $\{\pi_n\}$ に対する学習アルゴリズムを求める。正の推移確率行列 q の全体 $Q = \{q = (q_{ij}^k) | q_{ij}^k > 0, \sum_{j=1}^L q_{ij}^k = 1\}$ をパラメータ

空間として、 q の推定量として $\text{MLE } q(n) = (q_{ij}^k(n))$ を使う。ここで、 $q_{ij}^k(n) = n_{ij}^k / \sum_{j \in S} n_{ij}^k$ 、ただし n_{ij}^k は n 期の履歴 $h_n = (i_0, a_0, \dots, i_n)$ において $i_t = i, a_t = k, i_{t+1} = j$ となる t の個数。初期情報 $q(0) \in Q$ と $\lambda_n \in (0, 1)$ なる数列 $\{\lambda_n\}$ を適当に選んで n 期の情報様式を $\bar{q}(n) = \lambda_n q(0) + (1 - \lambda_n) q(n)$ で表わす。各 $k \in A$ と $q \in Q$ に対して、 L 次元ユークリッド空間 R^L の作用素 $U^k[q] = (U_1^k[q], \dots, U_L^k[q])$ を

$$(4.5) \quad U_i^k[q] u = r(i, k) + \sum_{j \in S} (q_{ij}^k - \eta(q)) u_j, \\ u = (u_1, \dots, u_L) \in R^L$$

$$\text{ただし } \eta(q) = \min_{i, j \in S, k \in K} q_{ij}^k$$

で定める。このとき $\eta(q) > 0$ であるから、作用素 $U[q] = \text{may}_{k \in A} U^k[q]$ は縮小写像となる。今、不動点を $u^* = (u_1^*, u_2^*, \dots, u_L^*)$ とすれば、 $u^* = U[q] u^*$ が成り立ち、この式において $\eta^* = \eta(q) \sum_{j \in S} u_j^*$ とおけば、平均最適基準のもとでの最適方程式を得る：

$$(4.6) \quad u_i^* = \max_{k \in A} \{r(i, k) - \eta^* + \sum_{j \in S} q_{ij}^k u_j^*\} \quad (i \in S)$$

各 $i \in S$ に対して、上式の右辺を最大にする $k \in A$ の全体を $A_i^*[q]$ で表わすと $f(i) \in A_i^*[q]$ ($i \in S$) となる任意の定常政策 f は q に対して定まる MDP の平均最適となり、 η^* はそのときの最適平均期待利得となる ([5, 6, 18])。今、推定量 $\bar{q}(n)$ を使って L 次元ベクトルの列 $\{\tilde{V}(n) = (\tilde{V}_1(n), \tilde{V}_2(n), \dots, \tilde{V}_L(n))\}$ を逐次決めてゆくことによる (value-iteration)。

$$(4.7) \quad \tilde{V}_i(0) = 0$$

$$\tilde{V}_i(n+1) = \max_{k \in A} U_i^k[\bar{q}(n)] \tilde{V}(n) \quad (n \geq 0)$$

次に $b_0 = 1, b_n > b_{n+1} > 0$ ($n \geq 0$) なる数列 $\{b_n\}$ に対して、 $\phi(b_n) = b_{n+1}$ ($n \geq 0$) を満たす増加関数

$\phi : [0, 1] \rightarrow [0, 1]$ によって, π_{n-1} を π_n に次のように改訂する (学習アルゴリズム). (4.7) の右辺を最大にする任意の k を $\tilde{k}_{n+1}(i)$ で表わす. 各 $i \in S$ に対して, $\tilde{k}_{n+1}(i) = k_i$ ならば

$$(4.8) \quad \pi_n(k_i|i) = 1 - \sum_{l \neq k_i} \phi(\pi_{n-1}(l|i))$$

$$\pi_n(l|i) = \phi(\pi_{n-1}(l|i)) \quad (l \neq k_i).$$

このアルゴリズムでは, $\pi_n(k_i|i) > \pi_{n-1}(k_i|i)$, $\pi_n(l|i) < \pi_{n-1}(l|i)$ ($l \neq k_i$) が成り立つので(4.8)は一種の reward-penalty 型の学習方式となっている. (4.7) と (4.8) から $\pi_0(\cdot|i)$ を適当に与えれば 1 つの適応政策 $\pi = (\pi_0, \pi_1, \dots)$ が構成される. このとき, 任意の $i \in S$, $k \in A$ に対して, $\pi_0(k|i) > 0$ かつ $\lim_{n \rightarrow \infty} \lambda_n = 0$, $\lim_{n \rightarrow \infty} b_n = 0$, $\sum_0^\infty b_n = \infty$ が成り立つならば, 推定量 $\bar{q}(n)$ の一致性と確率 1 で $\lim_{n \rightarrow \infty} \bar{V}_n = u^*$, $\lim \sum_{k \in A^*|q_i} \pi_n(k|i) = 1$ が示される. これは, π_n はいくらでもパラメータの値 q に対する平均最適な政策に近づくことを意味している. また π^* は平均最適な適応政策でもある ([10]).

5. おわりに

適応型の問題では, 未知のパラメータに関する情報の収集・処理およびシステムの最適化という 2 つの側面があり, 両者を統一的にとり扱うのが適応制御過程といえる. いわば, 統計学の諸理論と OR の決定モデルの解析手法とが融合した新しい理論といえよう. したがって, 適応制御的な手法はもっと現実の社会の中に浸透していてもおかしくないと思えるのだが意外にそうではない. これは, たとえば問題を DP で定式化しても, 一般に状態空間が高次元の集合になり実際に解を求めるのが困難のためであろう. しかし, 計算アルゴリズムの開発と高速電算機の利用により, 適応型の問題解決の手法は現実の問題処理に対して有効な戦力になりつつあると思われる.

参考文献

[1] Bellman, R. E.: *Dynamic Programming*. Princeton University Press, 1957

[2] Bellman, R. E.: *Adaptive Control Processes; A Guided Tour*. Princeton Univ. 1962

[3] Federgruen, A. and Schweitzer, P. T.: Non-stationary Markov decision problems with converging parameters. *J.O.T.A.* 34 (1981)

[4] 古川・門田: マルコフ決定過程の展望, 第 4 回数理解計画シンポジウム論文集, 1983, 111-141

[5] ハワード, R・A 著, 関根他訳: ダイナミックプログラミングとマルコフ過程, 培風館, 1971

[6] 金子哲夫: マルコフ決定理論入門, 槇書店

[7] Kurano, M.: Discrete-time MDP with an Unknown Parameter. *J. Oper. Res. Soci. Japan*, 15 (1972), 67-76

[8] Kurano, M.: Adaptive Policies in MDP with Uncertain Transition Matrices. *J. Inf. & Opti. Sci.* 4 (1983), 21-40

[9] 蔵野・安田・中神: 不確実情報の MDP と応用, 第 4 回数理解計画シンポジウム論文集, 1983, 159-178

[10] Kurano, M.: *Learning Algorithms for MDP (Preprint)*, 1985

[11] Lakshmivarahan, S.: *Learning Algorithms: Theory and Applications*. Springer, 1981

[12] Mandl, P.: Estimation and Control in Markov Chains. *Adv. Appl. Prob.* 6 (1974), 40-60

[13] 宮沢光一: 情報・決定理論序説, 岩波, 1971

[14] 大野勝久: マルコフ過程の計算アルゴリズム, 第 4 回数理解計画シンポジウム論文集, 1983

[15] Petkau, A. T.: Sequential Medical Trials for Comparing an Experimental with a Standard Treatment. *J.A.S.A.* 73 (1978), 328-338

[16] 佐藤他: ベイズ的手法を用いた未知パラメータを含むマルコフ決定過程の漸近的性質. 電気通信学会論文誌, J 61-D, 1978, 1-8

[17] 坂口 実: 動的計画法, 至文堂, 1968

[18] 坂本武司: マルコフ決定過程, 情報科学講座「マルコフ過程」共立出版, 1966, 106-169

[19] Van Hee, K. M.: *Bayesian Control of Markov Chains*, Mathe. Centre Tracks, Amsterdam, 1978

[20] Yakowitz, S. J.: *Mathematics of Adaptive Control Processes*. Elsevier, New York, 1969