

マルコフ・システムの制御

大野 勝久

1. はじめに

ORにおける代表的な確率システムといえば、在庫管理、待ち行列、信頼性であり、これらはすべてマルコフあるいはその一般化としてのセミ・マルコフ過程として論じられてきた。経済、工学等の他の分野においても、確率の変動が無視できない動的な問題をあつかうさいには、マルコフ過程に帰着させて論ずるのが通例である。このような過程で記述されるマルコフ・システムの制御問題の例として、多品目在庫管理問題を考えよう。

例1. 多品目在庫管理問題

問題は、期のはじめに各品目の在庫量をみて、発注するかしないか、発注するとすればどれだけかを与える最適発注政策を求めることである。費用としては、発注費用、品切れ費用、在庫費用を考え、発注時にかかる固定費用は注文する品目の組合せに依存するものとする。したがって、必ずしも (σ, S) 政策が最適とはならず、混合発注政策等を考える必要がある[4, 7, 8]。さて、品目数を c で表わし、 n 期における品目 $l(=1, \dots, c)$ の在庫量を x_l^n 、需要量を z_l^n 、発注量を a_l^n で表わそう。発注品は次期のはじめまでに納品されるものとするれば、 $n+1$ 期のはじめにおける在庫量 x_l^{n+1} は、

$$x_l^{n+1} = \max\{x_l^n - z_l^n, 0\} + a_l^n \quad (1)$$

で与えられる。需要量 (z_1^n, \dots, z_c^n) が各期独立に

分布するものと仮定すれば、 x_l^{n+1} の確率分布は x_l^n, a_l^n と z_l^n の周辺分布から計算できる。同様に、需要分布から n 期にかかる平均総費用が計算できる。このように現在の状態 (x_1^n, \dots, x_c^n) と決定 (a_1^n, \dots, a_c^n) が与えられれば、その期にかかる平均総費用と次期の状態 $(x_1^{n+1}, \dots, x_c^{n+1})$ の確率分布が求められる。

今1つの例として、図1に示される $M/M/c$ 待ち行列システムの制御問題を考えよう[12]。

例2. $M/M/c$ 待ち行列システムの制御

客は到着率 λ のポアソン過程にしたがって到着し、各窓口でのサービス時間は指数分布にしたがって、待ち行列長は L に制限されている。問題はシステムの状態に応じて、各窓口 $l(=1, \dots, c)$ のサービス率を M_l+1 個のレベル $\{0, \mu_l, \dots, M_l \mu_l\}$ から選び、さらに複数の窓口が遊休中(サービス率が正で空)であればどの窓口に客を割り当てるかを与える最適制御政策を求めることである。客の待ち費用、窓口の稼働費用、遊休費用、サービス

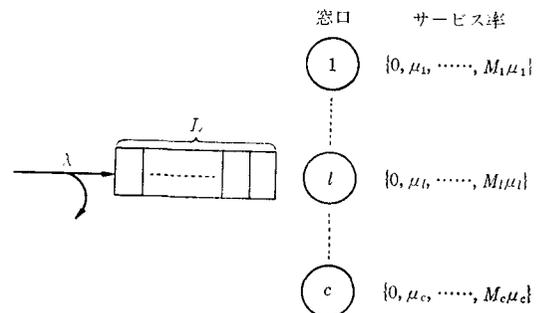


図1 $M/M/c$ 待ち行列システム

率の切り換え費用を考え、状態として待ち行列長、各窓口の状態（空かどうか）と現在のサービス率^{註1}をとる。決定は各窓口のサービス率および客の割り当てである。このとき例1と同様、現在の状態と決定が与えられれば、次に遷移する状態の分布とその遷移までにかかる平均総費用を計算することができる。しかし、この問題における決定はサービス率であり、(1)式の発注量のように直接的に状態を制御できず、サービス分布のパラメータとして確率的に制御できるにすぎない。このことが在庫管理問題との基本的な相違である。

例1、例2を含むマルコフ・システムの制御問題を統一的に論ずる理論が（セミ・）マルコフ決定過程である。以下では実際に最適制御政策を計算するアルゴリズムに焦点を絞り、簡単のためマルコフ決定過程について述べる。セミ・マルコフ決定過程への拡張は有限期間問題を除けば直接的である。また、例1、2とも費用最小化問題であったが、大勢にしたがい利得最大化問題を取りあつかうことにする。

2. 有限期間問題

システムは M 個の状態 $i \in I = \{1, \dots, M\}$ をとり、各状態 i においてとりうる決定は集合 $A_i = \{1, \dots, K_i\}$ で与えられている。このとき例1、2におけるように、状態 i で決定 $k \in A_i$ をとったとき次期の状態が j となる確率 $p_{ij}(k)$ およびその期における平均利得 $r_i(k)$ が与えられる。 $p_{ij}(k) \geq 0$, $\sum_{j \in I} p_{ij}(k) = 1$ である。単位期間当りの利子率 r から $\beta = 1/(1+r)$ として決められる $\beta (0 < \beta \leq 1)$ を割引率と呼ぶ。 n 期のはじめにおける状態を x_n , その時の決定を a_n で表わせれば、 x_n, a_n が与えられたときの n 期における平均利得 y_n は、

$$y_n = r_x(a_n), \text{ ここで } x = x_n$$

となり、その現価換算値は $\beta^n y_n$ である。問題は初期状態 $x_0 = i$ が与えられたとき、 N 期までに得られる総利得

$$\sum_{n=0}^N \beta^n y_n$$

の平均を最大にする決定列 (a_0, \dots, a_N) を定める最適政策を導くことである。

アルゴリズム^{註2}

1: すべての $i \in I$ に対して、

$$v_i^N = \max_{k \in A_i} r_i(k) \quad (2)$$

を計算し、最大を与える k を f_i^N とおく。 $n = N-1$ とおく。

2: $v_i^{n+1} (i \in I)$ の値を用い、各 $i \in I$ に対して、

$$v_i^n = \max_{k \in A_i} \{r_i(k) + \beta \sum_{j \in I} p_{ij}(k) v_j^{n+1}\} \quad (3)$$

を計算し、最大を与える k を f_i^n とおく。

3: $n = n-1$ とおき、 $n < 0$ となれば停止。さもなければステップ2へもどる。

v_i^0 が最大平均利得を与え、 $f^n = (f_1^n, \dots, f_M^n)$ とおけば、求める最適政策は (f^0, f^1, \dots, f^N) で与えられる。すなわち、期間0では決定 f_i^0 をとり、 n 期 ($= 1, \dots, N$) のはじめに状態が j となれば決定 f_j^n をとるのが最適である。この政策は、システムの履歴 (x_0, a_0, \dots, x_n) を情報として用い、決定を確率的に定める非常に広い政策のなかで最適である。このアルゴリズムの計算量は、

$$\text{乗算: } MN(M+1), \text{ 加算: } M^2 N,$$

$$\text{大小比較: } (N+1) \left\{ \sum_{i=1}^M K_i - M \right\} \quad (4)$$

であり、必要とするコア・メモリーは $\{v_i^n\}, \{v_i^{n+1}\}$ の記憶に要する $2M$ である。また最適政策を1つ求めればよいのであれば、最適政策の記憶に $(N+1)M$ のメモリーが必要である。したがって、たとえ連続状態をとるシステムであっても、最適政策の性質がくわしくわかっていないかぎり、それを計算する段階で必然的に有限状態で近似せざるを得ない。

上記アルゴリズムの計算量は N に比例して増えるだけであるが、得られた最適政策を実際用いるためには、 f^0, \dots, f^N の表を必要とし、期間 n が変わるとに対応する f^n の表を参照して決定を決めなければならない。したがって、計画期間長 N が長い問題では、あたかも永久にシステムをと

りまく現在の状況が維持されるものと考え、 $N \rightarrow \infty$ として得られる最適定常政策を採用するのが実際のであろう。

3. 割引利得問題

計画期間長が無限であり、割引率 $\beta < 1$ のとき得られる総利得

$$\sum_{n=0}^{\infty} \beta^n y_n$$

の平均を最大にする最適定常政策 (f^*, f^*, f^*, \dots) を求める問題である。最適定常政策のもとでは、どの期間 n においても $x_n = i$ となれば決定 f_i^* をとればよく、ただ1つの表 f^* を必要とするだけである。この問題に対しては政策反復法、線形計画法、逐次近似法が古典的な結果として知られている[2]。

政策反復法

1: 初期政策 $f^0 = (f_1^0, \dots, f_M^0)$ を与え、 $n=0$ とおく。

2: (値決定ルーチン) M 元連立一次方程式

$$v_i - \beta \sum_{j \in I} p_{ij}(f_i^n) v_j = r_i(f_i^n) \quad (i \in I) \quad (5)$$

を解き、 $v_i (i \in I)$ を求める。

3: (政策改良ルーチン) 各 $i \in I$ に対して

$$\max_{k \in A_i} \{r_i(k) + \beta \sum_{j \in I} p_{ij}(k) v_j\} \quad (6)$$

を与える決定の集合 F_i^{n+1} を求める。もし $f_i^n \in F_i^{n+1}$ ならば $f_i^{n+1} = f_i^n$ とおき、さもなければ f_i^{n+1} を F_i^{n+1} の任意の要素にとる。すべての i で $f_i^{n+1} = f_i^n$ となれば停止。 f^n は最適定常政策 f^* であり、 $v_i (i \in I)$ は状態 i から出発したときの最大割引利得 v_i^* を与える。ある i で $f_i^{n+1} \neq f_i^n$ となれば、 $n=n+1$ としてステップ2へもどる。

政策反復法は有限回の反復で f^* を与えることが知られている。このアルゴリズムが N_1 回の反復で停止したときの計算量は、ステップ2でガウスの消去法を用いるとすれば、^{注3)}

$$\text{乗除算: } N_1 \left\{ \frac{1}{3} M^3 + O(M \sum_{i=1}^M K_i) \right\} \quad (7)$$

$$\text{加減算: 同上, 比較: } N_1 \left\{ \sum_{i=1}^M K_i - M \right\}$$

となり、必要とするコア・メモリーは $M^2 + 4M$ である。

線形計画法

$\sum_{i \in I} b_i = 1$ を満たす適当な正数 b_i を与え、次のLP問題:

$$\begin{aligned} \max \quad & \sum_{i \in I} \sum_{k \in A_i} r_i(k) x_{ik} \quad (8) \\ \text{subject to} \quad & \sum_{k \in A_i} x_{ik} - \beta \sum_{j \in I} \sum_{k \in A_j} p_{ji}(k) x_{jk} = b_i \\ & x_{ik} \geq 0 \quad (i \in I, k \in A_i) \end{aligned}$$

を解く。最適解 x_{ik}^* は各 i に対して唯一の $k \in A_i$ で正となり、その k を f_i^* ととれば最適定常政策 f^* が得られる。

LP問題(8)は $\sum_{i=1}^M K_i$ 変数 M 制約式問題である。また、政策反復法の初期政策 f_i^0 に対応する x_{ik} を基底変数としてとれば、初期基底実行可能解が得られる。このとき改訂単体法を用いて N_2 回で最適解が得られたものとすれば、計算量は⁴⁾、

$$\text{乗除算: } M^3 + O(N_2 M \sum_{i=1}^M K_i) \quad (9)$$

$$\text{加減算: 同上, 比較: } N_2 \left\{ \sum_{i=1}^M K_i - M \right\}$$

となり、基底逆行列をコア・メモリーに保持すれば $M^2 + O(M)$ が必要となる。線形計画法では反復当たり1つの状態の決定しか改良されず、 N_2 は N_1 にくらべて相当大きく、 M のオーダーになることを注意しなければならない。

ここで、例1において品目数 $c=3$ ととり、各品目の在庫の上限を9ととれば、 $M=10^3$, $\sum_{i=1}^M K_i = 10^5$ となる。ゆえに政策反復法では千元連立一次方程式を解かねばならず、計算量は(7)式より $(N_1/3) \cdot 10^9$ のオーダーとなる。一方線形計画法では、10万変数千制約式のLP問題となり、計算量は(9)式より $10^9 + O(N_2 \cdot 10^8)$ のオーダーとなる。このように、いずれのアルゴリズムを用いるにせよ、状態数 M が数百以上となれば実際に最適政策を計算するのは容易ではない。

連立一次方程式の数値解法では、数百元以上の問題に対してヤコビ法、ガウス・ザイデル法等の反復解法がよく用いられている。したがって、こ

れら反復解法を値決定ルーチン (ステップ2) で用いる政策反復法が考えられる。しかし, ステップ3では改良された政策が得られる程度の(5)の近似解で十分であり, 反復解法を有限回の反復 m で打ち切ったほうが効率的であろう。すなわち, ステップ2で反復解法を m 回用いる政策反復法が得られる。この方法を修正政策反復法と呼び, 特に $m=0$ のとき古典的な逐次近似法が得られる。以下のアルゴリズムでは, 最大割引利得の上, 下限を計算し, それらを用いて最適になり得ない決定を A_i から除去する非最適政策の除去法を併用している。

修正政策反復法 [11]

1: 初期値 $v_i^0 (i \in I)$, 非負数 ε , 非負整数 m を与え, $A_i^0 = A_i, l_i = -\infty (i \in I), \xi = 0, n = 0$ とおく。

2: 各 i に対して,

$$w_i^{n+1} = \max_{k \in A_i^n} \{r_i(k) + \beta \sum_{j \in I} p_{ij}(k) v_j^n\} \quad (10)$$

を計算し, 同時に最大を与える決定の集合 F_i^{n+1} と

$$A_i^{n+1} = \{k \in A_i^n; r_i(k) + \beta \sum_{j \in I} p_{ij}(k) v_j^n \geq l_i - \beta \xi\}$$

を定める。もし $f_i^n \in F_i^{n+1}$ ならば $f_i^{n+1} = f_i^n$ とおき, さもなければ f_i^{n+1} を F_i^{n+1} の任意の要素にとる。すべての i で F_i^{n+1} が唯一の要素 f_i^{n+1} になればステップ6へ。

3: $z_i^0 = w_i^{n+1}$ とおき, $l = 0, 1, \dots, m-1$ に対して

$$z_i^{l+1} = r_i(f_i^{n+1}) + \beta \sum_{j \in I} p_{ij}(f_i^{n+1}) z_j^l \quad (i \in I) \quad (11)$$

を計算し, $v_i^{n+1} = z_i^m (i \in I)$ とおく。

4: $D = \max_{i \in I} \{v_i^{n+1} - v_i^n\}, \varphi = \min_{i \in I} \{v_i^{n+1} - v_i^n\}$

$$a = \max_{i \in I} \{v_i^{n+1} - w_i^{n+1}\}, b = \min_{i \in I} \{v_i^{n+1} - w_i^{n+1}\} \\ \xi = ((\beta D - b) / (1 - \beta)) \quad (12)$$

$$\eta = \max\{\beta \varphi - a\} / (1 - \beta), b \beta^m / (1 - \beta^m) \quad (13)$$

を計算し, $\xi - \eta < 2\varepsilon$ となればステップ6へ。さもなければ, $l_i = v_i^{n+1} + \eta (i \in I), n = n+1$ とおいてステップ2へもどる。

5: 最大割引利得 v_i^* の ε 近似値が,

$$v_i = v_i^{n+1} + \frac{1}{2}(\xi + \eta) \quad (i \in I) \quad (14)$$

で与えられ, f^{n+1} は ε' 最適政策である。ここで,

$$\delta = \max_{i \in I} \{v_i - r_i(f_i^{n+1}) - \beta \sum_{j \in I} p_{ij}(f_i^{n+1}) v_j\} \quad (15)$$

$$\varepsilon' = \varepsilon + \{\delta / (1 - \beta)\}. \quad (16)$$

6: f^{n+1} は唯一の最適定常政策である。

$$c = \max_{i \in I} \{w_i^{n+1} - v_i^n\}, d = \min_{i \in I} \{w_i^{n+1} - v_i^n\}$$

を計算し, $\varepsilon'' = \beta(c - d) / (1 - \beta)$ とおけば,

$$v_i = w_i^{n+1} + \{\beta(c + d) / 2(1 - \beta)\} (i \in I) \quad (17)$$

は最大割引利得 v_i^* の ε'' 近似値を与える。

このアルゴリズムは任意の $v_i^0 (i \in I)$ と $\varepsilon > 0$ に対して有限回の反復で停止し, ε' 最適政策 (ステップ5) あるいは唯一の最適政策 f^* (ステップ6) を与える。特に最適政策が唯一であることがあらかじめ知られていれば, $\varepsilon = 0$ ととることで有限回の反復で f^* を求めることができる。修正政策反復法が N_3 回の反復で停止したときの計算量は, $|A_i^n|$ で A_i^n の要素数を表わすことにすれば注5)

$$\text{乗算} : (M+1) \left\{ \sum_{n=1}^{N_3} \sum_{i=1}^M |A_i^{n-1}| + mN_3 \right\} + o(M)$$

$$\text{加減算} : M \left\{ \sum_{n=1}^{N_3} \sum_{i=1}^M |A_i^{n-1}| + (m+3)N_3 \right\} + o(M)$$

$$\text{比較} : 2 \sum_{n=1}^{N_3} \sum_{i=1}^M |A_i^{n-1}| + 3MN_3 + o(M) \quad (18)$$

となる。必要なコア・メモリーは v_i^n 等に $4M, A_i^n$ のために $\sum_{i=1}^M K_i$ ビット必要である。たとえば例1で $M=10^3, \sum_{i=1}^M K_i=10^5$ の場合, 計算量は反復の初期では 10^8 であるが終期では 10^3 のオーダーにまで低下し有力な手法であることが示される。

例2で $c=2, L=10, M_1=M_2=4$ の場合 $M=740, \sum_{i=1}^M K_i=4980$ となるが, FACOM M-200による計算結果では, $m=6(12)$ のとき25(15)回の反復で唯一の最適政策が16(13)秒で得られた。なお, 修正政策反復法では(10), (11)式でガウス・ザイデル法の反復を用いることもでき, $m=6(12)$ のとき13(9)回の反復で9(8)秒に短縮された。 m を反復ごとに自動的に調整する修正政策反復法が, 例2と同様な直列型待ち行列システムの制御問題に応用されている[9]。

4. 時間平均利得問題

問題は無限期間における単位時間当り利得

$$\lim_{N \rightarrow \infty} \frac{1}{N+1} \sum_{n=0}^N y_n$$

の平均 (ゲインと呼ばれる) を最大にする最適常政策を求めることである。

政策反復法

1: 初期政策 f^0 を与え, $n=0$ とおく。

2: (値決定ルーチン) $g_i^n, v_i^n (i \in I)$ に関する次の

$2M$ 元連立一次方程式を解く。

$$g_i^n - \sum_{j \in I} p_{ij}(f_i^n) g_j^n = 0 \quad (i \in I) \quad (19)$$

$$v_i^n - \sum_{j \in I} p_{ij}(f_i^n) v_j^n = -g_i^n + r_i(f_i^n) \quad (i \in I) \quad (20)$$

3: (政策改良ルーチン) 各 i に対して,

$$G_i^{n+1} = \left\{ \sum_{j \in I} p_{ij}(k) g_j^n \text{ を最大化する } k \in A_i^n \right\}$$

$$F_i^{n+1} = \left\{ r_i(k) + \sum_{j \in I} p_{ij}(k) v_j^n \text{ を最大化する } k \in G_i^{n+1} \right\}$$

を定める。 $f_i^n \in F_i^{n+1}$ ならば $f_i^{n+1} = f_i^n$ とおき, さもなければ f_i^{n+1} を F_i^{n+1} の適当な要素にとる。すべての i で $f_i^{n+1} = f_i^n$ となれば停止し, そうでなければ $n=n+1$ とおいてステップ 2 へもどる。

$P(f)$ で定常政策 f から定まる遷移確率 $p_{ij}(f_i)$ を i 行 j 列要素とする行列を表わせば, $P(f)$ はマルコフ連鎖を与える。この連鎖のエルゴード (再帰) 集合の個数を $e(f)$ とおくと, 行列 $(I-P(f))$ の階数は $M-e(f)$ となる^{註6)}。ゆえに (19), (20) で解 g_i^n は一意に決まるものの, v_i^n は $e(f^n)$ 個の自由度をもつ。Howard[5] は各エルゴード集合に属する 1 つの状態 で $v_i^n = 0$ とおき, 得られた $2M-e(f^n)$ 元連立一次方程式を解けばよいとしている。しかし, そのためには反復ごとにマルコフ連鎖の状態分類をしなければならず, アルゴリズム[3] は知られているものの決して効率的ではない。この点を改善したものが次の値決定アルゴリズムであり,

$$\Delta g_i = g_i^n - g_i^{n-1}, \Delta v_i = v_i^n - v_i^{n-1} \quad (i \in I)$$

に関する連立一次方程式を解いている。

値決定アルゴリズム[14]

1: 行列 $(I-P(f^n))$ にガウスの消去法等を用い, 適当に状態を並べ換えて,

$$L(I-P(f^n)) = \begin{pmatrix} Q_1 & Q_2 \\ & 0 \end{pmatrix} \quad (21)$$

なると $M \times M$ 下三角行列 L , $(M-e(f^n)) \times (M-e(f^n))$ 上三角行列 Q_1 , $(M-e(f^n)) \times e(f^n)$ 行列 Q_2 を求める。以下添字 1 で Q_1 に対応する $M-e(f^n)$ 次元ベクトルを表わし, 添字 2 で残りの $e(f^n)$ 次元ベクトルを表わすことにする。

2: $c_1 = (P(f^n)g^{n-1})_1 - g_1^{n-1}$ (22)

$$d_l = r(f^n)_l + (P(f^n)v^{n-1})_l - g_l^{n-1} - v_l^{n-1} \quad (l=1, 2) \quad (23)$$

を計算し, $c_1 = d_2 = 0$ となればステップ 4 へ。

3: $Q_1 \Delta g_1 + Q_2 \Delta g_2 = c_1$, $(L \Delta g)_2 = d_2$ (24)

を解き, その解 Δg (M 次元ベクトル) を用いて,

$$Q_1 \Delta v_1 = d_1 - (L \Delta g)_1 \quad (25)$$

を解く。 $\Delta v_2 = 0$ とおきステップ 6 へ。

4: (23) 式の d_1 で要素が正となる状態の集合を S とおき, 次の集合^{註7)}

$$T = \{i \in I; \text{ある } n \geq 0 \text{ で } \sum_{j \in S} p_{ij}^n(f^n) > 0\}$$

を最短経路を求める Dijkstra のアルゴリズム (たとえば [1]) を用いて求める。

5: $i \in T$ に対して $\Delta v_i = 0$ とおいて,

$$Q_1 \Delta v_1 + Q_2 \Delta v_2 = d_1 \quad (26)$$

を解き, $\Delta v_i (i \in T)$ を求める。 $\Delta g_i = 0 (i \in I)$ とおく。

6: $g_i^n = g_i^{n-1} + \Delta g_i$, $v_i^n = v_i^{n-1} + \Delta v_i (i \in I)$ (27) とおく。

このアルゴリズムを政策反復法の値決定ルーチンで用いれば有限回の反復で最適定常政策 $f^* = f^n$ および最大ゲイン $g^* = g^n$ が得られる。

線形計画法[6]

$\sum_{i \in I} b_i = 1$ を満たす適当な正数 $b_i (i \in I)$ を与え, 次の LP 問題:

$$\max \sum_{i \in I} \sum_{k \in A_i} r_i(k) x_{ik} \quad (28)$$

$$\text{subject to } \sum_{k \in A_i} x_{ik} - \sum_{j \in I} \sum_{k \in A_j} p_{ji}(k) x_{jk} = 0 \quad (i \in I)$$

$$\sum_{k \in A_i} (x_{ik} + y_{ik}) - \sum_{j \in I} \sum_{k \in A_j} p_{ji}(k) y_{jk} = b_i \quad (i \in I)$$

$$x_{ik} \geq 0, y_{ik} \geq 0 \quad (i \in I, k \in A_i)$$

を解き最適解 x_{ik}^* , y_{ik}^* を求める. 最適政策 f_i^* は,

$$(i) \sum_{k \in A_i} x_{ik}^* > 0 \text{ ならば } x_{ik}^* > 0 \text{ となる任意の } k$$

(ii) $\sum_{k \in A_i} x_{ik}^* = 0$ ならば $y_{ik}^* > 0$ となる任意の k として得られる.

政策反復法では M 元連立一次方程式 (24) を解かねばならず, 線形計画法では $\sum_{i \in I} K_i$ 変数 $2M$ 制約式の LP 問題を解かなければならない. したがって, M が大きな問題では逐次近似法によらざるを得ないものと思われる.

逐次近似法 [13]

1: 初期値 $v_i^0 (i \in I)$, 正数 ϵ および

$$0 < \tau < \min\{1/(1-p_{ii}(k)); p_{ii}(k) < 1, i \in I, k \in A_i\}$$

を満たす τ を適当に与える. $n=0$ とおく.

2: 各 i に対して,

$$g_i^{n+1} = \max_{k \in A_i} \{r_i(k) + \sum_{j \in I} p_{ij}(k) v_j^n - v_i^n\} \quad (29)$$

を求め, $\Delta = \max_{i \in I} g_i^{n+1}$, $\nabla = \min_{i \in I} g_i^{n+1}$ を計算し, $\Delta - \nabla < \epsilon$ となればステップ 4へ.

3: $v_i^{n+1} = v_i^n + \tau(g_i^{n+1} - g_i^n)$ ($i \leq M-1$) (30)

$$v_M^{n+1} = 0$$

とおき, $n=n+1$ としてステップ 2へもどる.

4: (29) 式で最大を与える k を f_i とおけば, $f = (f_1, \dots, f_M)$ は最大ゲイン g^* との差が ϵ 以下となる ϵ 最適政策を与える.

このアルゴリズムの収束を保証するには, 残念ながら, 最大ゲイン g^* が状態に依存しない定数であることを仮定しなければならない. この仮定は f^* のもとでのマルコフ連鎖がただ 1つのエルゴード集合をもてば成立し, 例 1 では常に, 例 2 ではサービス率の切り換え費用が低ければ成立する.

時間平均利得問題に対する修正政策反復法は現在のところ公表されていないが, 筆者は前節に述べた導出と同じ考え方で, $m=0$ で逐次近似法と一致する修正政策反復法を得ている. この反復法は $e(f^*) \geq 2$ となる一般的な問題へも拡張できる.

5. おわりに

有限状態, 有限決定をもつマルコフ決定過程におけるアルゴリズムについて, 現在までに得られた結果を紹介した. 紙数の関係もあり, 引用すべき文献を多数省略してしまったのは心残りである. 最後に, この解説では状態が完全に観測できるものとしてきたが, 実際の問題では観測が不完全な場合も多く, この場合については最近の解説 [10] を参照されたい.

引用文献

- [1] エイホ, A. V. 他(野崎, 野下訳): アルゴリズムの設計と解析, サイエンス社, 1977
- [2] Derman, C.: *Finite State Markovian Decision Processes*. Academic Press, New York, 1970
- [3] Fox, B. L. and Landi, D. M.: An Algorithm for Identifying the Ergodic Subchains and Transient States of a Stochastic Matrix. *Comm. ACM*, 11(1968), 619-621
- [4] Goswick, T. E. and Sivazlian, B. D.: The Mixed Ordering Policy in Periodic Review Stochastic Multi-Commodity Inventory Systems. *Naval Res. Log. Quart.*, 21(1974), 389-410
- [5] ハワード, R. A. (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] Hordijk, A. and Kallenberg, L. C. M.: Linear Programming and Markov Decision Chains. *Management Sci.*, 25(1979), 352-362
- [7] Johnson, E. L.: Optimality and Computation of (σ, S) Policies, in the Multi-Item Infinite Horizon Inventory Problem. *Management Sci.*, 13(1967), 475-491
- [8] Kalin, D.: On the Optimality of (σ, S) Policies, *Math. Operations Res.*, 5(1980), 293-307
- [9] 三根, 大野, 市木: 直列型待ち行列システムの最適制御に関する計算アルゴリズムについて, OR学

会春季アブストラクト集. 175-176, 1982年3月

- [10] Monahan, G. E. : A Survey of Partially Observable Markov Decision Processes : Theory, Models, and Algorithms. *Management Sci.*, 28(1982), 1-16
- [11] Ohno, K. : A Unified Approach to Algorithms with a Suboptimality Test in Discounted Semi-Markov Decision Processes. *J. Operations Res. Soc. Japan*, 24(1981), 296-324
- [12] 大野勝久: 待ち行列システムの最適制御に関する計算アルゴリズムについて, 数理解析研究所講究録 452, 1-19, 1982年2月
- [13] Platzman, L. : Improved Conditions for Convergence in Undiscounted Markov Renewal Programming. *Operations Res.*, 25(1977), 529-533
- [14] Spreen, D. : A Further Anticycling Rule

in Multichain Policy Iteration for Undiscounted Markov Renewal Programs, *Zeitschrift für Operations Res.*, 25(1981), 225-233

注)

- 1) サービス率の切り換え費用の計算に必要である.
- 2) $p_{ij}(k)$, $r_i(k)$ が n に依存する非定常な場合にも適用できる.
- 3) $O(M)$ は $M \rightarrow \infty$ のとき M と同位の無限大を表わす.
- 4) 理論的には $O(M^{2.81})$ のアルゴリズムがあるが, 実用的ではない[1].
- 5) $O(M)$ は $M \rightarrow \infty$ のとき M の低位の無限大を表わす.
- 6) I は単位行列を表わす.
- 7) $p_{ij}^n(f^n)$ は $P(f^n)^n$ の i 行 j 列要素を表わす.

特集に当って

大野 勝久

「確率システム? 知らないナァ」とページをめくり「ウン」と頭をかかえられた読者もおありかと思いますが, 現実に解決をせまられている問題の多くは, 多かれ少なかれはっきりしない不確定部分を含んでおります. このような問題を状態あるいは決定変数を用い, 不確定部分を確率として定式化したものが確率システムです (少々オーバーないいかたかも知れませんが, また定式化されたからといって現在の理論, 計算機で必ず解けるわけでもありませんが……). 自動制御の分野では, 「確率システム・シンポジウム」(日本自動制御協会主催)が毎秋, 専門の異なった参加者が集まって開かれており, 昨秋で13回と回を重ねております. そこで, このシンポジウムをはじめられた砂原先生に, その理論的側面を概説していただき, 中

おおの かつひさ 京都大学

溝, 片山両先生にメイン・テーマである安定性と推定をお願いいたしました.

一方, ORをながめてみますと, 待ち行列, 信頼性, 在庫管理等確率システムとよぶことのできる各々永い歴史をもった分野があります. 研究の進展につれ, 専門化, 細分化するのは学問の常ですが, 現在では相互の交流, 批判もなく孤立化しているようにみうけられます. そこでこれら分野で横断的に用いられている (セミ・マルコフ過程をとりあげ, 高橋先生に「マルコフ分析」をお願いいたしました).

なお, 同様な趣旨でOR学会関西支部に「応用確率論研究部会」(主査, 西田俊夫先生)が設けられ, 今春から活動しております. ぜひご参加ください.

おわりに, ご多忙のところをご執筆いただいた上記諸先生に厚くお礼申しあげます.