

氏名(本籍)	小野謙二(大分県)
学位の種類	工学博士
学位記番号	工博第219号
学位授与年月日	昭和45年3月25日
学位授与の要件	学位規則第5条第1項該当
研究科専門課程	東北大学大学院工学研究科 (博士課程)電気及通信工学専攻
学位論文題目	マルコフ決定過程に関する研究

(主査)
 論文審査委員 教授 本多 波雄 教授 福島 弘毅
 教授 大泉 充郎 助教授 木村 正行

論文内容要旨

1 まえがき

近年、自動制御理論の1つの分野である確率論的(stochastic)制御理論は急速に発達し、多くの新しい制御理論や制御方式が研究されている。しかし、これらの理論を実用化の立場から眺めてみると、応用の実例は甚だしく少ない。その1つの理由として、制御対象における未知のパラメータ(例えば確率分布の平均値、分散など)の推定と確率論的制御理論が別々の立場で論じられてきたことがあげられる。

本論文の目的は、未知パラメータの推定と確率論的制御理論を有機的に結合する新しい制御方式を開発することである。すなわち、観測により得られたデータから未知パラメータの推定値を逐次求めていく段階においても、システムの評価関数(損失とする)が最小となるように逐次制御し

ていく方式を与えた。

考察の対象とするシステムとしては、離散的で有限な状態がマルコフ連鎖をなすシステムの制御、すなわちマルコフ決定過程を論じた。この問題は大別して、推移確率行列が既知である場合と未知である場合に分類される。また、状態を観測する立場から、真の状態を観測できる場合と雑音などが付加され真の状態が観測できない（状態観測が不完全な）場合に別けられる。本論文では、離散型有限マルコフ決定過程において、状態の推移確率行列は既知であるが、状態観測が不完全な場合、および状態の推移確率行列が未知の場合（真の状態が観測可能と不能なとき）について考察した。

2 推移確率行列が既知の決定過程

定常、非周期的、既約なマルコフ連鎖の推移確率行列を

$$P(k) = [p_{ij}^k] \quad (1)$$

$$\sum_{j=1}^C p_{ij}^k = 1 \quad i = 1, 2, \dots, C \\ k = 1, 2, \dots, A$$

であらわす。ここに p_{ij}^k は、状態 x_i のとき制御行動 u_k をとることにより状態 x_j に推移する確率をあらわす。また、状態と観測値の対応をあらわす確率行列を

$$Q = [q_{ij}] \quad (2)$$

$$\sum_{j=1}^D q_{ij} = 1 \quad i = 1, 2, \dots, C$$

とあらわす。ここに q_{ij} は、状態が x_i のとき観測値 y_j が得られる確率である。このように、推移確率行列が既知で、観測値が状態の確率的関数であらわされる場合、状態と観測値の対を 1 つの状態を考えて新しいマルコフ連鎖を構成する。この推移確率は

$$P'(k) = [p_{(i,j)(i',j')}^k] \quad (3)$$

$$p_{(i,j)(i',j')}^k = p_{ii'}^k q_{i'j'}$$

となる。以上の結果から、状態 (x_i, y_i) で制御行動 u_k を決定したとき状態 $(x_{i'}, y_{j'})$ への推移にともなう損失を $r_{(i,j)(i',j')}^k$ と定義する。この損失と新しく定義された推移確率行列 $P'(k)$ に、ダイナミック・プログラミング法を適用した解法を示した。

3 推移確率行列が未知の決定過程

真の状態は観測できるが、状態の推移確率行列が未知な場合として、次のような制御を考えることができる。すなわち、状態 x と観測値 y の対を事象 z と定義する。この事象の系列が r 重マルコ

フ連鎖をなし， r が既知であるとき確率的近似法を適用して直接にシステムの評価関数 θ の期待値を計算し逐次制御する。また r が未知の場合の制御方法，及び電子計算機によるシミュレーション結果を示した。

3.1 評価関数の期待値

システムの評価関数 θ の期待値の推定値は，確率的近似法を適用して

$$\begin{aligned}\hat{E}_{n(i,k)}[\theta | \eta_i, u_k] &= \hat{E}_{n(i,b)-1}[\theta | \eta_i, u_k] \\ &+ r_{n(i,k)} \{ \theta(t+1) - \hat{E}_{n(i,k)-1}[\theta | \eta_i, u_k] \} \\ r_{n(i,k)} &= 1/n(i,k)\end{aligned}\quad (4)$$

ここで $n(i,k)$ は，時刻 $t+1$ までに事象の系列 η_i が観測され，かつ制御行動 u_k を行なった回数である。また，

$$\theta(t+1) = g(\eta_i, u_k, y_j) \quad (5)$$

である。ここで y_j は時刻 $t+1$ における観測値をあらわす。(4)式は， θ の平均値と分散が有限であれば $n(i,k) \rightarrow \infty$ で期待値に確率収束することが証明される。

3.2 制御行動の決定

長さ r なる任意の事象の系列 η_i に対して，システムの評価関数 θ の期待値を最小とする制御行動 u を求めたい。すなわち，

$$E[\theta(t+1) | \eta_i, u] = \min_{u'} \{ E[\theta(t+1) | \eta_i, u'] \} \quad (6)$$

を満足する最適な制御 u を求める。このような最適制御行動 u を決定する関数は，任意の事象の系列 η_i に対する条件付確率 $P_r[u | \eta_i]$ で与えられる。この確率の推定値は， θ の期待値の推定値 $\hat{E}_n[\theta | \eta_i, u]$ を利用し，次式で逐次修正し推定できる。

$$\hat{P}_{n(i)}[u_k | \eta_i] = \alpha \hat{P}_{n(i)-1}[u_k | \eta_i] + (1-\alpha) \psi_{n(i)}(u_k, \eta_i) \quad (7)$$

$$\psi_{n(i)}(u_k, \eta_i) = \begin{cases} 1 & \hat{E}_{n(i,k)}[\theta | \eta_i, u_k] \\ & = \max_{u_1} \{ \hat{E}_{n(i,k)}[\theta | \eta_i, u_1] \} \\ 0 & \text{他の場合} \end{cases}$$

ここで $n(i)$ は，時刻 $t+1$ までに事象の系列 η_i が観測された回数をあらわす。 α は $0 < \alpha < 1$ なる定数とする。上式は $n(i) \rightarrow \infty$ で最適な制御行動を決定する確率 $P_r[u_k | \eta_i]$ に確率収束することが証明される。

3.3 制御手順

事象の系列が r 重マルコフ過程をなす場合の制御手順を示す。

(1) 初期セット

任意の制御 u_k と事象の系列 η_i に対して、初期値をそれぞれ

$$\hat{P}_o[u_k | \eta_i] = 1/A \quad (8)$$

$$\hat{E}_o[\theta | \eta_i, u_k] = 0 \quad (9)$$

とする。ここで A は制御行動の数。

(2) 制御の決定

事象の系列（長さ r ）を η_i とすると、時刻 $t+1$ における制御行動 u_k を確率 $\hat{P}_{n(i)}[u_k | \eta_i]$ にしたがってランダムに決定する。

(3) 評価関数 θ ($t+1$) の計算

時刻 $t+1$ における制御行動 u_k と観測値 y_j と事象の系列 η_i に対する評価関数を(5)式によつて計算する。

(4) θ の条件付期待値

θ の条件付期待値の推定値を(4)式によって求める。

(5) 最適な制御行動の決定

事象の系列 η_i について、全ての制御行動に対する θ の期待値の推定値を比較して $\psi_{n(i)}$ (u, η_i) を決める。

(6) 決定関数の修正

全ての制御行動 u に対して、行動の決定確率を(7)式により修正する。

(7) ステップを進める。

事象の系列を新しくし(2)へとぶ。

3.4 多重マルコフ過程の制御

これまで事象の系列が r 重マルコフ連鎖をなすと仮定したが、この r が未知である場合について考察した。この場合、事象の長さ 0 から始めて、マルコフ性の検定結果にもとづき r を逐次大きくしていく方法を考えた。また検定として、確率密度をもっているということの他は何も規定しない連(r un)の数を用いるのが有効である。このような r が未知のシステムに対して、電子計算機を利用し簡単なシミュレーションを行なった。

4 マルコフ連鎖の関数

状態の集合を $\Omega_x = \{x_1, x_2, \dots, x_c\}$ 、観測値の集合を $\Omega_y = \{y_1, y_2, \dots, y_D\}$ 、
 f を Ω_x から Ω_y への関数とする。ここで $C \geqq D$ とする。また任意の観測値 v に対応する $f^{-1}(v)$

で表わされる状態の集合を考え、その要素の数を N_ν 、任意の長さの観測値の系列を v_1, v_2, \dots, v_n , w_1, w_2, \dots, w_m であらわす。観測値の系列 $v_i \cdot \nu \cdot w_j$ が得られる確率を $p(v_i \nu w_j)$ であらわす。次に、

$$P_\nu(v_1, v_2, \dots, v_n; w_1, w_2, \dots, w_m) = \prod_{i=1}^n p(v_i \nu w_j) \quad (10)$$

$$j = 1, 2, \dots, m$$

が正則であるような最大の整数 n を n_ν とすると、次の関係がある。

- (1) 任意の観測値 ν に対して $n_\nu \leq N_\nu \leq C$
- (2) $C \geq D$ ならば $\sum_\nu n_\nu \leq C$
- (3) $\sum_\nu n_\nu = C \Leftrightarrow n_\nu = N_\nu$

確率的関数の場合、(10)式が正則であるような系列 $v_1, v_2, \dots, v_c, w_1, w_2, \dots, w_c$ の集合のなかで長さが最少な系列を $v_1^*, v_2^*, \dots, v_c^*, w_1^*, w_2^*, \dots, w_c^*$ とすると

- (4) $\sum_\nu n_\nu \leq C \times D$
- (5) $\max_i \{ \lg(v_i^*) \} \leq C$

$$\max_j \{ \lg(w_j^*) \} \leq C$$

ここで $\lg(v)$ は、系列 v の長さを示す。

観測値の系列から、状態の推移確率行列を同定することは一般に不可能である。また、状態の関数としてあらわされる観測値の系列の分布から、次の観測値を推定することも一般には不可能である。

5 む す び

マルコフ決定過程において、推移確率行列の既知と未知、状態観測が完全と不完全な場合を論じた。本論文の決定過程は、その構造的一般性により、最適制御のみならず、ORなど応用分野も広いものである。

審査結果の要旨

マルコフ決定過程の理論は、確率的制御のみならず、広く情報科学、計画理論、オペレーションズ・リサーチなどの分野においても数学的基礎として欠くことのできないものである。しかし、実用化の立場から眺めてみると、この理論を応用した例は甚だしく少ない。その主な理由として、制御の対象となる系の未知パラメータの確率分布や推移確率行列などの推定と制御理論を関連づける研究がなかったことがあげられる。

本論文は、系の状態の推移がマルコフ連鎖で記述できる場合の制御、すなわち、マルコフ決定過程について、未知パラメータの推定と制御理論とを有機的に結合した新しい制御方式について、著者が研究した成果をまとめたもので全編5章よりなる。

第1章は緒論であって、本研究の目的とともに、マルコフ決定過程についての従来の研究と本研究との関連が述べられている。

第2章では、推移確率行列が既知で観測値が状態の確率的関数として表わされる場合について、状態とその観測値を対とする新しいマルコフ過程を定義することにより、ダイナミック・プログラミングの手法を用いて最適決定を行なう方法を与えていた。

第3章では、推移確率行列は未知であるが状態の観測が完全な場合のマルコフ過程の制御について新しい制御方式を与えていた。この制御方式は、観測可能な量と制御可能な量の対で表わされる事象が n 重マルコフ過程をなす一般的な場合を扱っており、確率的近似法により直接には推移確率行列を推定することなく評価関数を逐次最小にする制御を行なうもので、推定と制御とを有機的に結合した優れた制御方式であり、本研究の重要な成果である。また、本章では「連」によるマルコフ連鎖の次数 n を検定する方法を与え、これを上記の制御方式に結合して、その実用性を一層高めている。なお、シミュレーションにより本方式の有効性が確かめられている。

第4章では、推移確率行列が未知で状態観測が不完全な系の制御を行なう手がかりとして、マルコフ連鎖の関数について2, 3の有用な性質を導いている。

第5章は結論である。

以上要するに、本論文は確率的近似法と「連」によるマルコフ性の検定法を導入して、統計的推定の理論と制御理論とを有機的に結合することにより、高度な制御機能をもつ実用性の高い制御方式を開発したもので、その成果はシステム工学および制御工学の発展に寄与するところが少なくない。

よって、本論文は工学博士の学位論文として合格と認める。