

## 統計的学習理論の応用

### に関する研究

阪大 工 魚崎 勝司 今村 秀樹  
田坂 誠男 杉山 博

#### §1 はじめに

近年における制御理論の発展の一つの側面は いわゆる適応制御, 学習制御に関する理論の開拓であらうと思われた。

すなわち 従前の理論では 'given' ('known') としてきた制御対象の特性, 構造, また入力, 外乱の特性などを本来の姿である 'unknown' なものとしてとらえて, analysis や synthesis をおこなうという立場にたつての理論の開拓と発展であり, これは記述函数的表現や非線形理論の発展とも関連している。

このとき重要なものは "学習" 過程というもののとりえ方である。"学習" は 明らかに過去の行動(こゝには環境, 刺激, 応答, 結果なども含む)と密接に結びついたものとして 行動の変化を示す現象であり, この観点から "学習" 過程を構成することが必要である。

我々はここへ、三年にわたって、心理学において学習行動の解析の手段として提案されていた R. Bush と F. Mosteller とによる統計的学習理論を“学習”の手法として応用することとを考へ、理論的、実験的な考察をおしすすめてきた。ここではその一端を紹介したい。

### 3.2 統計的学習理論の概要<sup>1)</sup>

人間や動物の学習行動を説明するための試みは、従来から心理学的な立場から、種々行なわれてきている。以下で述べる R. Bush と F. Mosteller による統計的学習理論もその一つである。彼らは W. Estes による刺激標本モデル (Stimulus sampling model) を整理、拡張して学習オペレータなるものを導入し、学習過程を一つの確率過程におおきかえて解析をすすめている。彼らの学習過程のモデルを図示すると Fig. 1 のようになる。

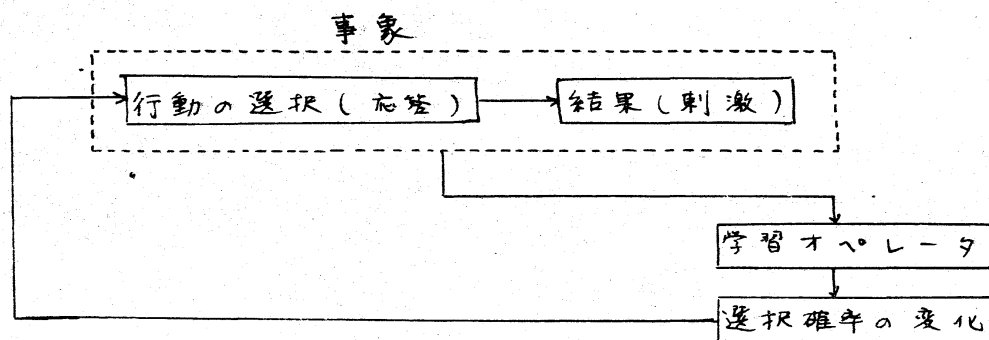


Fig. 1. Bush-Mosteller 学習モデル概念図

行動の選択  $A$  と結果  $O$  とを組み合わせた一つの事象  $E$  に対応して、選択確率を変化させる学習オペレータがあるのがあるが、その形は一般には、 $r$  個の選択と、 $s$  個の結果の種類があるときには、

$$(2.1) \quad Q_{ij} p = \alpha_{ij} p + a_{ij}, \quad i=1, 2, \dots, r, \quad j=1, 2, \dots, s$$

または

$$(2.2) \quad Q_{ij} p = \alpha_{ij} p + (1 - \alpha_{ij}) \lambda_{ij}, \quad \begin{matrix} i=1, 2, \dots, r, \\ j=1, 2, \dots, s \end{matrix}$$

こゝに、学習オペレータの作用を受けたのも確率としての性質を維持するにため、

$$0 < \alpha_{ij} \leq 1, \quad 0 \leq \lambda_{ij} \leq 1 \quad \begin{matrix} i=1, 2, \dots, r \\ j=1, 2, \dots, s \end{matrix}$$

の制限が、パラメータに付加される。

(2.1) からわかるように、このオペレータは形式的には、一般的なオペレータ

$$Qp = a_0 + a_1 p + a_2 p^2 + \dots$$

の近似をなしている。

これらのオペレータは、選択、結果によって確率的に作用されることになり、その確率を  $p_{ij}$  と表わすとき、この  $p_{ij}$  のとり方により、学習に3つの型があることが示される。以下では  $r=2$ ,  $s=2$  の場合について述べるが、その他の場合も同様である。

(i) Experimenter controlled events

ある事象の起る確率が、実験者によって定められていて、  
 試行によっても変化しない場合で、 $p_{ij}$ はその時点の選択の確  
 率  $p$  に依存しない。学習オパレーターは

$$(2.3) \quad \begin{aligned} Q_1 p &= \alpha_1 p + (1 - \alpha_1) \lambda_1 && \text{with probability } \pi_1 \\ Q_2 p &= \alpha_2 p + (1 - \alpha_2) \lambda_2 && \text{with probability } \pi_2 = 1 - \pi_1 \end{aligned}$$

となる。この場合の試行のくり返しによる選択確率の変化の  
 ようすを Fig. 2 に示す。

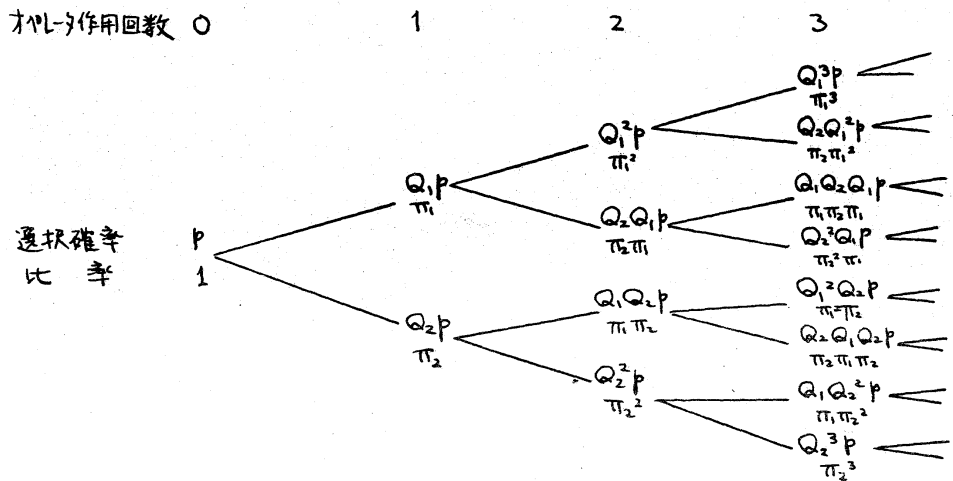


Fig. 2 確率の推移

(ii) Subject-controlled events

選択と事象とが一致する場合で、 $p_{ij}$ はその時点の選択の確  
 率  $p$  そのものになる。学習オパレーターは次で与えられる。

$$(2.4) \quad \begin{aligned} Q_1 p &= \alpha_1 p + (1 - \alpha_1) \lambda_1 && \text{with probability } p \\ Q_2 p &= \alpha_2 p + (1 - \alpha_2) \lambda_2 && \text{with probability } q = 1 - p \end{aligned}$$

この場合の選択確率の変化のようすを Fig. 3 に示してある。

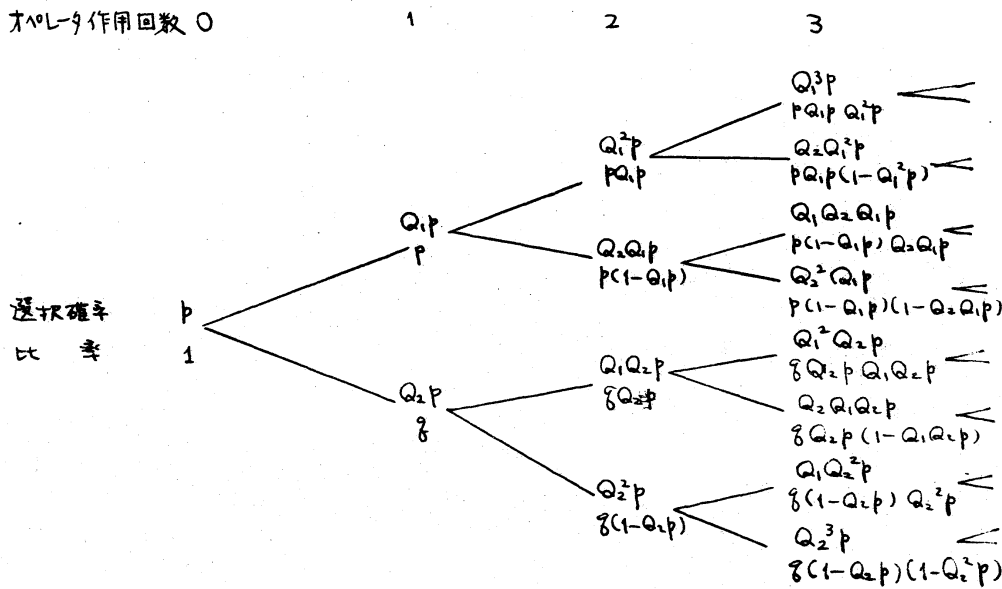


Fig. 3 確率の推移

(iii) Experimenter - Subject - controlled events

事象の起る確率が一定であり (i) と、事象  $\alpha$  起る確率  $\pi$ 、選択の確率に等しい (ii) の二つの考え方を結びつけたもので、ある選択がなされたとき、その選択がなされたという条件のもとで、ある結果の起る条件付確率が一定でありとする。学習オパレータは

$$\begin{aligned}
 Q_{11} p &= \alpha_{11} p + (1 - \alpha_{11}) \lambda_{11} && \text{with pr. } p \pi_1 \\
 Q_{12} p &= \alpha_{12} p + (1 - \alpha_{12}) \lambda_{12} && \text{with pr. } p (1 - \pi_1) \\
 Q_{21} p &= \alpha_{21} p + (1 - \alpha_{21}) \lambda_{21} && \text{with pr. } \delta \pi_2 = (1 - p) \pi_2 \\
 Q_{22} p &= \alpha_{22} p + (1 - \alpha_{22}) \lambda_{22} && \text{with pr. } \delta (1 - \pi_2) = (1 - p) (1 - \pi_2)
 \end{aligned}$$

(2.5)

によって与えられる。

Bush と Mosteller は, これらのモデルにおいて適当なパラメータを決めることにより, いくつかの学習実験のデータに対しよい一致を得ることができると示し, データからパラメータを統計的に推定するいくつかの方法についても若干の考察をおこなっている。

### §3 proportional parameter の推定に対する応用

二つの事象 A, B が比率  $p$ ,  $q (= 1-p)$  で起る系列  $x$  と与えられたとき, A の生起の比率のパラメータ  $p$  を推定するための手段として, Experimenter controlled events の場合の Bush-Mosteller learning model を採用した場合について考察した結果を述べる。

(i)  $\lambda_1 = 1$ ,  $\lambda_2 = 0$  とし, ある先験的知識より定める初期値  $\hat{p}_0$  ( $0 \leq \hat{p}_0 \leq 1$ ) より出発し, A, B の生起にしたがって次のオペレータ  $Q_{(1)}$ ,  $Q_{(2)}$  を作用させる。

$$(3.1) \quad \begin{aligned} \hat{p}_{n+1} &= Q_{(1)} \hat{p}_n = \alpha \hat{p}_n + 1 - \alpha, & \text{if A occurs} \\ \hat{p}_{n+1} &= Q_{(2)} \hat{p}_n = \alpha \hat{p}_n, & \text{if B occurs} \end{aligned}$$

このとき得られる系列  $\{\hat{p}_n\}$  は

$$(3.2) \quad \hat{p}_n = \sum_{i=1}^n (1-\alpha) \alpha^{i-1} y_i + \alpha^n \hat{p}_0, \quad n=1, 2, \dots$$

こゝに  $\{y_i\}$  は

$$(3.3) \quad \begin{aligned} y_i &= 1 & \text{if A occurs} \\ y_i &= 0 & \text{if B occurs} \end{aligned}$$

を確率変数列として表わされるので、

$$(3.4) \quad \begin{aligned} E(\hat{p}_n) &= (1-d^n)p + d^n \hat{p}_0 \\ &\rightarrow p \quad n \rightarrow \infty \end{aligned}$$

が言える。しかし

$$(3.5) \quad \begin{aligned} \text{Var}(\hat{p}_n) &= p(1-p)(1-d^{2n})(1-d)/(1+d) \\ &\rightarrow p(1-p)(1-d)/(1+d) \quad n \rightarrow \infty \end{aligned}$$

で、一般に、分散は  $n \rightarrow \infty$  としても 0 とはならないので、学習の一つの評価函数である真値からのずれ  $E(\hat{p}_n - p)^2$  も  $n \rightarrow \infty$  で 0 とはならない。

(ii) (3.1) のオパレータは Bush - Mosteller が考えたよりも、広い意味にとり、 $d$  が time-dependent, すなわち  $n$  に依存するモデルとして採用すると、

$$(3.6) \quad \begin{aligned} \hat{p}_{n+1} &= Q_{1|n} \hat{p}_n = \alpha_n \hat{p}_n + (1-\alpha_n) \quad \text{if } A \text{ occurs.} \\ \hat{p}_{n+1} &= Q_{2|n} \hat{p}_n = \alpha_n \hat{p}_n \quad \text{if } B \text{ occurs.} \end{aligned}$$

の学習オパレータとなる。ここで種々の criterion に従って、適切なパラメータ系列  $\{\alpha_n\}$  を選定することは重要である。

(3.6) の学習オパレータを用いると、 $\{\hat{p}_n\}$  は

$$(3.7) \quad \hat{p}_n = \sum_{i=1}^n (1-\alpha_i) \prod_{j=i+1}^n \alpha_j y_i + \prod_{i=1}^n \alpha_i \hat{p}_0 \quad n=1, 2, \dots$$

となる。よって

$$(3.8) \quad E(\hat{p}_n) = (1 - \prod_{i=1}^n \alpha_i) p + \prod_{i=1}^n \alpha_i \hat{p}_0$$

$$(3.9) \quad \text{Var}(\hat{p}_n) = p(1-p) \left\{ (1-\alpha_n)^2 + \sum_{i=1}^{n-1} (1-\alpha_i)^2 \prod_{j=i+1}^n \alpha_j^2 \right\}.$$

したがって  $n$  までには少なくとも 1 個の整数  $k$  が存在し、

$$(3.10) \quad \alpha_k = 0$$

ならば、

$$(3.11) \quad E(\hat{p}_n) = p$$

として、不偏性が示される。

この不偏性の条件下で最小分散を与えるパラメータ列は、

$$(3.12) \quad \alpha_n = 1 - (1/n)$$

なることが導かれる。このパラメータを用いた学習では、次式の成立するようになる。

$$(3.13) \quad \begin{aligned} E(\hat{p}_n) &= p \\ \text{Var}(\hat{p}_n) &= E(\hat{p}_n - p)^2 = p(1-p)/n \end{aligned}$$

これは、 $p$  の最尤推定量のもう性質と一致している。

(iii) 単に  $E(\hat{p}_n - p)^2$  を最小にするパラメータ列は、

$$(3.14) \quad \alpha_n = 1 - (1/(n + \hat{a})),$$

ただし

$$\hat{a} = p(1-p) / (p - \hat{p}_0)^2$$

により与えられる。この学習では、

$$(3.15) \quad \begin{aligned} E(\hat{p}_n) &= p / (1 + \hat{a}/n) + \hat{p}_0 \hat{a} / (n + \hat{a}) \\ \text{Var}(\hat{p}_n) &= np(1-p) / (n + \hat{a})^2 \\ E(\hat{p}_n - p)^2 &= p(1-p) / (n + \hat{a}) \end{aligned}$$

である。



(iv) 真値  $p$  は学習の前段階では未知だからパラメータ列として、上述の  $\{\alpha_n\}$  をとることができないので、この近似的として、パラメータ列  $\{\alpha_n\}$  を次のようにとることを考える。

$$(3.16) \quad \alpha_n = 1 - (1/(n+a)) \quad n=1, 2, \dots$$

この学習では、初期値  $\hat{p}_0$  の選択について

$$(3.17) \quad p - \{2p(1-p)/a\}^{1/2} \leq \hat{p}_0 \leq p + \{2p(1-p)/a\}^{1/2}$$

の関係が成立していれば、 $E(\hat{p}_n - p)^2$  は (ii) の学習より小となる。

(v) また学習理論における regression 関係 (3.8) を用いて、過去の  $\{\hat{p}_i\}$  の系列から、 $p$  を最小二乗推定することは可能で、この学習による推定は、

$$(3.18) \quad \hat{p}_n = \frac{\sum_{i=1}^{n-1} (1 - \prod_{j=1}^i \alpha_j) (\hat{p}_{i+1} - (\prod_{j=1}^i \alpha_j) \hat{p}_0)}{\sum_{i=1}^{n-1} (1 - \prod_{j=1}^i \alpha_j)^2}$$

であり、次の性質をもっている。

$$(3.19) \quad E(\hat{p}_n) = p$$

$$\text{Var}(\hat{p}_n) = E(\hat{p}_n - p)^2 = p(1-p) \sum_{i=1}^n (1 - \alpha_i)^2 \left\{ \sum_{k=1}^i (1 - \prod_{l=1}^k \alpha_l) \prod_{l=k+1}^i \alpha_l \right\}^2 / \left\{ \sum_{j=1}^i (1 - \prod_{l=1}^l \alpha_l)^2 \right\}^2$$

なお最初の先験的分布がベータ分布  $B(\bar{a}, \bar{b})$  としたときの Bayesian type の学習過程を、この問題に適用すると、推定値は、

$$\bar{p} = \int_0^1 p^{m+1} (1-p)^{n-m} dH(p) / \int_0^1 p^m (1-p)^{n-m} dH(p)$$

$$= B(m+\bar{a}+1, n-m+\bar{b}) / B(m+\bar{a}, n-m+\bar{b})$$

となり、(iv) において

$$a = \bar{a} + \bar{b}$$

とした場合に相当することには注意する。

以上の各種の学習と比較すると

不偏性のもとで分散小 (ii)  $\succ$  (v)

$E(\hat{p}_n - p)^2$  小 (iii)  $\succ$  (iv)  $\succ$  (ii)  $\succ$  (v)

この議論は、三個以上の事象の生起確率の推定にも拡張できる。

#### § 4 Markovian sequence の予測問題への応用

$r$  個の状態をもつ単純 Markov chain の遷移確率行列は

$$(4.1) \quad P = \begin{bmatrix} p_{11} & p_{12} & \cdots & p_{1r} \\ p_{21} & p_{22} & \cdots & p_{2r} \\ \vdots & \vdots & \ddots & \vdots \\ p_{r1} & p_{r2} & \cdots & p_{rr} \end{bmatrix}$$

で与えられるが、これを学習オペレータを用いて推定するには、時点  $n$  で状態  $E_i$  にあって、時点  $n+1$  で状態  $E_j$  に遷移したなら、遷移確率行列の第  $i$  行の成分  $p_{ik}$  の推定値  $\hat{p}_{ik}$  を、学習オペレータにより、次のように変化せしめる。

$$(4.2) \quad \begin{aligned} \hat{p}_{ij} &\rightarrow \alpha \hat{p}_{ij} + 1 - \alpha \\ \hat{p}_{ik} &\rightarrow \alpha \hat{p}_{ik} \quad k \neq j, \quad k=1, \dots, r \end{aligned}$$

ただし  $\alpha$  は前節で述べたような意味にとり、このとき、前節の議論を若干変更するだけで、成分が row-wise に推定されること加えられる。

この手法を、 $s$  個の事象  $A_1, \dots, A_s$  を含む、周期  $r$  の系列 (



遷移確率行列の推定のようにすと Fig. 4 に、10 回毎の誤り個数の推移を Fig. 5 に示す。この結果は 10 回のくり返し実験の平均である。ただ問題となるのは、周期が長く事象の種類が多いときで、このときには、状態数が非常に大きくなることによる難点が生じてくる。

これから pure strategy の場合には、 $\alpha$  による影響はそれほど大きくないように見える。しかし別の解析から、 $\alpha$  を小さくすれば、分散が大きくなり、誤り危険率が高くなることが示されるので、 $\alpha$  は 1 に十分近くとることを望ましいと結論される。これに反し、mixed strategy では全体的に pure strategy より誤り率が高いことが明らかで、これは予期されることである。また学習の初期では、 $\alpha$  の小さなほど誤る割合は小さいが、これは  $\alpha$  のもう意味からもうなづける結果である。

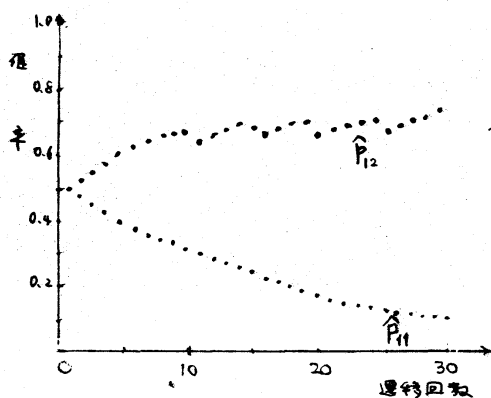


Fig. 4.  $P$  の第 1 行成分の推定  
のようす  $\alpha=0.95$

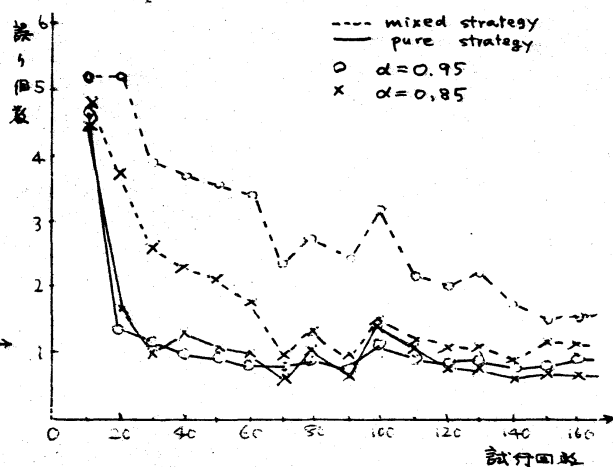


Fig. 5. 誤り回数の変化

### § 5 Two-armed bandit problem の応用

Two-armed bandit problem とは、二種類の選択の可能性(以

下では2個の coin, AとBを想定する。それぞれが、その各々の選択において、1 unit の gain を得る(表が出る)確率  $\pi_1, \pi_2$  が定まっているが、実験者には未知であるとするとき、 $N$ 回の試行で、gain をできるだけ大きくするには、どのように選択していけば良いかということを取り扱う問題で、主に Bayesian の立場から考察がなされている。ここでは、min-max approach, Bayesian approach が採用されて考察がなされているが、これは状況が不明のときの制御問題の取り扱いと同様である。ここでは、自然との game と考えて、min-max approach を採用するが、特に、実際的な立場から、最初の何回かは、二つの選択の可能性を許すが、ある時点以降は、一方にのみきめて試行をおこなう場合の考察を行なう。

このような場合について、Vogel が一つの algorithm を提案している<sup>2)</sup>。それによると最初は coin A と coin B を対にしてふり、次の確率変数を導入する。

$$(5.1) \quad U_k = \sum_{i=1}^k (X_i - Y_i)$$

ここで  $X_i, Y_i$  はそれぞれ、coin A, coin B について、表が出たとき 1, 裏が出たとき 0 をとる確率変数である。さらにある正の整数  $\nu$  を導入し、

$$(5.2) \quad -\nu < U_k < +\nu$$

である限り、coin A と coin B の対の実験を行なうが、もし

$$(5.3) \quad U_k \geq +\nu \quad (\leq -\nu)$$

なら, 残り  $a$  ( $N-2k$ ) 回は coin A (coin B) のみを用いる。この

とき Loss function,

$$(5.4) \quad L(\sigma, \tau, \nu) = \sigma - SN/N \quad \sigma = \max(\pi_1, \pi_2), \tau = \min(\pi_1, \pi_2)$$

を考へ,  $\min_{\nu} \max_{\sigma, \tau} L(\sigma, \tau, \nu)$  を与える  $\nu$  とし

$$(5.5) \quad \nu = 0.292 N^{\frac{1}{2}} \quad N \geq 100$$

をとればよいことを示した。

我々は, これに対し, equal  $\alpha$  かつ  $\lambda_{11} = \lambda_{22} = 1$ ,  $\lambda_{12} = \lambda_{21} = 0$  の Experimenter-Subject controlled events の学習により, coin A を選択する確率を変化させることを考へた。すなわち

$$(5.6) \quad \begin{aligned} Q_{11}p &= \alpha p + 1 - \alpha && \text{with pr. } p\pi_1 \\ Q_{12}p &= \alpha p && \text{with pr. } p(1-\pi_1) \\ Q_{21}p &= \alpha p && \text{with pr. } (1-p)\pi_2 \\ Q_{22}p &= \alpha p + 1 - \alpha && \text{with pr. } (1-p)(1-\pi_2) \end{aligned}$$

のオペレータを用いる。ただし  $c$ ,  $1-c$  に吸収壁を設け,

$$(5.7) \quad p \geq 1-c \quad (\leq c)$$

ならばそれ以降は coin A (coin B) のみを用いることにした。こ

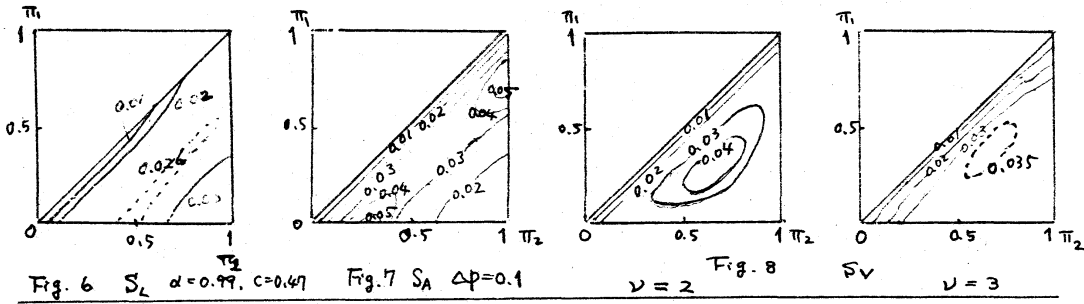
の方法 ( $S_2$ ) を用いたときの,  $N=100$  に対する Loss function の値を爆弾法により計算し, 図示したものが Fig. 6 である。なお,

coin A の選択確率を

$$(5.8) \quad \begin{aligned} Qp &= p + \Delta p && \text{with pr. } p\pi_1 + (1-p)(1-\pi_2) \\ &= p - \Delta p && \text{with pr. } p(1-\pi_1) + (1-p)\pi_2 \end{aligned}$$

で変化させた場合 ( $S_A$ ) については, Vogel の algorithm ( $S_V$ ) につ

いて Fig. 7, Fig. 8 に示した。



これから3つの strategy は  $\pi_1, \pi_2$  の組み合わせに対し、お互いに相補う性質をもつこととわかると共に、学習オロレータを用いた我々の方法は、Vogelの方法より、minmaxの意味でより良い性質を有していることが示された。

§6 結び

いくつかの例で示したように、統計的学習理論を用いた我々の方法は、望ましい学習に対する一つの指針を与えていることがわかる。ここでとりあげた例では、制御理論との結びつきが、それほど鮮明でないが、将来の応用には十分期待されるものがあると考ええる。たとえば、ここでとりあげた予測の問題は適応制御への一つの approach として、また Two-armed bandit problem への応用は identification の問題への一つの approach を示唆するものといえないであらうか。

今後の課題としては、ここでとりあげたいくつかの実験の理論的背景を追求するとともに、大局的な視野からみた制御理論への応用を考えていきたい。

## 参考文献

- 1) R. R. Bush, F. Mosteller, *Stochastic Models for Learning*. John Wiley  
(1955)
- 2) W. Vogel, "A sequential design for the two armed bandit," *Ann.  
Math. Statist.*, 31, 430 (1960)