

情報科学講座 A.2.5



情報理論 II

—情報の幾何学的理論—

北川敏男編

編集委員

大泉充郎
勝木保次
北川敏男
喜安善市
栗原俊彦
桑原万寿太郎
坂井利之
高田昇平
次田皓
南雲仁一
中村幸雄
和田弘

共立出版株式会社

1968

執筆者

甘利俊一 東京大学工学部

第5章 学習識別の理論

5.1 学習識別系

パターン認識系において、パターンの特徴信号の分布が知られている場合はむしろまれである。分布がわかっていなければ、前章で述べた方法で最適な決定方法を定めることはできない。また、たとえ分布がわかったとしても、それを用いて最適な決定方法を定めることは、実は容易でない。それに、パターンの分布は固定しているとは限らず、刻一刻と変動している場合も多い。このような状況を考えるとき、識別方法を計算で求めて固定しておくやり方は、あまり得策ではないことになる。これに代わるものとして、識別系に学習能力を付しておき、系にはいってくるパターン信号に基づいて学習を行なって、識別方法を逐次定めていく方法が考えられる。これが学習識別 (learning pattern-classification) の考えである。

学習識別には、おおざっぱに言って、次の二つの方法が考えられる。一つは過去のパターン信号のデータに基づいて、パターンの確率分布そのものを逐次求めていく方法であり、もう一つは、分布を媒介とせず、識別関数を直接求めていく方法である。

まず、前者について考えよう。過去に現われた、各クラスのパターン ξ をすべて記憶しておき、そのヒストグラムをつくれれば、確率分布 p_α や $p_\alpha(\xi)$ が近似的に求まる。 p_α や $p_\alpha(\xi)$ が求まれば、計算によって最適な識別ベクトルが求まることになる。しかし、過去に生じたパターンをすべて記憶しておくには、ばく大な容量の記憶装置が必要であるし、またその後の計算も大変なものである。そこで、通常は、確率分布の形を適当に仮定して、分布のパラメータを推定する方法が用いられる。

たとえば、分布 $p_\alpha(\xi)$ は正規分布であると仮定しよう。すると、分布 $p_\alpha(\xi)$

は C_α に属する ξ の平均値 $\bar{\xi}_\alpha$ と共分散行列 Σ_α とによって

$$p_\alpha(\xi) = \frac{1}{\sqrt{(2\pi)^n \det|\Sigma_\alpha|}} \exp\left\{-\frac{1}{2}(\xi - \bar{\xi}_\alpha)_t \Sigma_\alpha^{-1} (\xi - \bar{\xi}_\alpha)\right\} \quad (5.1)^*$$

と表わせる。それゆえ、過去のデータから正規分布のパラメータ $\bar{\xi}_\alpha$ と Σ_α を求めれば、分布が求まる。この方法は、分布のパラメータを推定することから、パラメトリックな方法とよばれる。

例として、 $\bar{\xi}_\alpha$ を求めよう。過去に C_α に属するパターンが k 回出たとし、それに基づいた $\bar{\xi}_\alpha$ の推定値を $\bar{\xi}_\alpha^k$ と書く。いままた、 C_α のパターン ξ が出たとすると、これを修正して

$$\bar{\xi}_\alpha^{k+1} = \frac{k\bar{\xi}_\alpha^k + \xi}{k+1} \quad (5.2)$$

とすれば、 $\bar{\xi}_\alpha$ の値が逐次求まっていくことになる。

分布自体を求めていくなれば、系についての正確な知識が得られようが、この方法には次のような欠点も考えられる。第一に、分布自体を推定するとなると、非常にたくさんのデータをたくわえねばならない。これを避けるためには、分布の形を事前に知らなければならない。第二に、分布がわかったとしても、それに基づいて識別ベクトルを求めることが容易でない。第三に、パターンの確率分布などの系の構造が変動しつつあるときには、この方法では変動に追従しにくいきらいがある。

したがって、この方法は系の構造が変動しない場合に、過去のデータより系に対する精密な知識を得たいときに用いられる。そして、それにより最適な識別ベクトルを求め、これを固定して使用する。系の構造が刻々と不規則に変化していく場合や、識別ベクトルを学習により直接求めたい場合は、この方法は不適當である。

そこで、分布を媒介とせず、識別ベクトルを学習により直接求める方法を考

* 行ベクトルと列ベクトルとを区別する必要があるときは ξ, θ は、列ベクトルを表わすものとし、添字 t は転置を意味する。

える。識別ベクトル θ を、現在系にはいつてくるパターン ξ とそれがどのクラスに属するかという情報に基づいて、逐次修正していく方法である。これは、分布の形を知る必要がなく、分布のパラメータの推定が不要であるため、**ノンパラメトリックな学習**といわれる。

ここでは、ノンパラメトリックな学習法を追求することとし、前者の分布の推定を含む学習法は、通常の統計学や他の学習識別の教科書⁴⁵⁾にゆずる。

ノンパラメトリックな方法として、いわゆる**パーセプトロンの学習法**⁴⁶⁾が有名である。しかし、これはパターンの分布が線形分離可能なときのみ収束する方法であり、線形識別関数を用いた場合に限定されている。ここでは、パターンの分布が重なっていてもよく、また一般の識別関数について成立する学習法である、**確率的降下法**について述べる。後節で、学習の収束、収束速度、収束精度、確率構造の変化に対する追従特性、などが明らかにされる。

確率的降下法の問題点として、次の点があげられる。まず、この方法が使えるのは、連続損失の場合に限られる。したがって、連続でない場合に用いたければ、損失を連続な関数で近似しなければならない*。第二に、確率的降下法による識別ベクトルは、 $L(\theta)$ の極小値に収束するが、これが必ずしも最小値とは限らない。したがって、最適な識別ベクトルを得るには、適当な初期ベクトルから出発しなければならないことがある**。

以上の点を考えると、この方法は、最適な識別ベクトルのだいたいの値はわかっているが、系の構造が変動しているために、細かい所までは定められないようなときに、特に有効といえる。音声認識装置は、特定の発声者に合わせて識別方法を固定したのでは、発声者が変わった場合にうまく働かない。識別装置に学習機能をもたせて、発声者の性質に合わせて、識別ベクトルを自動的に調整する場合などが、この例である。そのほか、プラントなどで、原料その他の状況の変動に応じて、制御方法を自動調整する場合なども、この例として考

* 近似を良くすると、必要な情報が境界 B_{gr} 上に局在するため、収束の速度が落ちる。4.3D 項参照。

** 線形分離可能なパターン分布に対しては、 $L(\theta)$ の極小値はただ一つであることが証明できる。

えられる。

5.2 確率的降下法による学習

A. 確率的降下法

識別ベクトル θ を、与えられたパターン ξ に基づいて逐次修正していく識別系を考える。新しく得られる識別ベクトルを θ' と書くと

$$\theta' = \theta + \delta\theta \quad (5.3)$$

であり、修正項 $\delta\theta$ は ξ および θ の属するクラスに依存する。

C_α のパターン ξ が提示されたときの修正項を

$$\delta\theta = \delta\theta(\xi, \alpha, \theta) \quad (5.4)$$

と書く。修正項を、 $L(\theta)$ を減少させるように選びたい。しかし、パターンの確率分布が未知であるために、 $L(\theta)$ も未知であり、 $L(\theta)$ を減少させる方向はわからない。そこで、次善の策として、 $L(\theta)$ を平均として減少させるように修正項を選ぶことを考える。「平均」とは、系にはいつてくるあらゆるパターン ξ についての未知の分布 $p_\alpha, p_\alpha(\xi)$ に基づいた期待値の意味である。

修正項の期待値は

$$\overline{\delta\theta} = \sum_\alpha \int p_\alpha p_\alpha(\xi) \delta\theta(\xi, \alpha, \theta) d\xi \quad (5.5)$$

と表わせる。 $\delta\theta$ は微小と考えると、1回の修正による $L(\theta)$ の変化は

$$\delta L(\theta) = \delta\theta \cdot \nabla L(\theta) \quad (5.6)$$

と書ける。 $\delta L(\theta)$ の期待値は

$$\overline{\delta L(\theta)} = \overline{\delta\theta} \cdot \nabla L(\theta)$$

であるから、これを負にするには、 ∇L と $\overline{\delta\theta}$ とが鈍角をなすようにする必要があり。このためには、 ϵ を正の定数、 C を正の定符号をもつ行列* として

$$\overline{\delta\theta} = -\epsilon C \nabla L(\theta) \quad (5.7)$$

* 任意のベクトル $a \neq 0$ に対して、常に $a_i C a_i > 0$ が成立する行列を、正の定符号をもつ行列という。

ならばよい。関数

$$\alpha_\alpha(\xi, \theta) = -\nabla L_\alpha(\xi, \theta) \quad (5.8)$$

の期待値は、(4.88) からわかるように

$$\overline{\alpha_\alpha(\xi, \theta)} = \sum_\alpha \int p_\alpha p_\alpha(\xi) \alpha_\alpha(\xi, \theta) d\xi = -\nabla L(\theta) \quad (5.9)$$

であるから*

$$\delta\theta(\xi, \alpha, \theta) = \varepsilon C \alpha_\alpha(\xi, \theta) \quad (5.10)$$

とおけば、(5.7) が成立する。 $\alpha_\alpha(\xi, \theta)$ を学習関数と名づけ、この学習法を確率的降下法という**。

この学習法では、より好ましい識別ベクトルが次々に得られる、というわけではない。修正の結果、識別ベクトルがまえよりも悪くなる場合も生ずる。しかし、(5.7) に示されるように、修正ベクトルの平均方向はより良くなる方向を向いている。

この方法を、 $L(\theta)$ の最小点を求めるのによく用いる、最急降下法と比較してみよう。最急降下法では、 $L(\theta)$ や $\nabla L(\theta)$ がわかっているために、 $L(\theta)$ を減少させる方向、すなわち山を下る方向がわかる。したがって、 $L(\theta)$ が必ず減るように修正を行なうことができる。これに反して学習問題においてはパターン分布が未知なために、 $L(\theta)$ や $\nabla L(\theta)$ がわからず、山を下る方向がわからない。しかし、識別装置にはいつてくるパターンは、未知のある確率分布に従って発生している。この情報を利用するために、パターンの関数 $\alpha_\alpha(\xi, \theta)$ で、その平均値が山を下る方向を向いているものをうまく探し出して、その方向へ確率的に下っていかうとするものである。

これは、坂の途中に立たされた「よっぱらい」が酔歩を行なう状況に似ている。彼は、あるときは山に登る方向によろめき、また次には山を下る方向によろける。登る確率よりは下る確率のほうが大きいため、彼は平均としては山を

* 本章では連続損失を考える。

** これは、確率的近似法 (stochastic approximation method) の応用と考えられる^{47,48)}。

ずり落ちて行き、ついには谷底へ落ち込む。そして谷底の近傍で、よろめき(微小振動)を続けるであろう。

確率的降下法により学習を行なわせると、同様の状況で、識別ベクトルが最適値の近傍に落ち着くことが予想される。収束の証明や、収束の精度(微小振動の状況)、収束の速度などは後に示される。

B. 種々の識別関数に対する学習法

a. 一次識別関数 初めに二分割の場合について述べる。この場合、識別関数は $g(\xi) = g_1(\xi) - g_2(\xi)$ を用いると

$$g(\xi, \theta) = \theta \cdot \hat{\xi}$$

と書いて、 g の正負に応じてパターンは C_1 もしくは C_2 と決定される。損失を

$$\left. \begin{aligned} l_1(\xi, \theta) &= l(-g) \\ l_2(\xi, \theta) &= l(g) \end{aligned} \right\} \quad (5.11)$$

とおこう*。学習関数は

$$\left. \begin{aligned} \alpha_1(\xi, \theta) &= -\nabla l_1 = l'(-g) \hat{\xi} \\ \alpha_2(\xi, \theta) &= -\nabla l_2 = -l'(g) \hat{\xi} \end{aligned} \right\} \quad (5.12)$$

となるから、次の学習規則が得られる。

$$\left. \begin{aligned} \delta\theta &= \varepsilon l'(-g) C \hat{\xi}, & \xi \in C_1 \text{ のパターンを識別したとき} \\ \delta\theta &= -\varepsilon l'(g) C \hat{\xi}, & \xi \in C_2 \text{ のパターンを識別したとき} \end{aligned} \right\}$$

特別な場合として、

$$l(d) = \begin{cases} d & d \geq 0 \\ 0 & d < 0 \end{cases} \quad C = E \text{ (単位行列)} \quad (5.13)$$

を考えると

$$\delta\theta = \pm \varepsilon \hat{\xi}, \quad \xi \text{ を誤識別したとき (+は } C_1, - \text{ は } C_2 \text{ のパターン) となる。}$$

* これは (4.89), (4.90), (4.91) で、 $g_1 \equiv g$, $g_2 \equiv 0$ とし、 $s_g = 1$ とおいたものである。

これは、パーセプトロン学習法として知られるものである⁴⁶⁾。

損失 (5.11) の幾何学的意味はあまり明確でない。損失として、識別面 B から誤識別されるパターン ξ までの距離 d を用いたほうが意味が明確である。パターン ξ から B までの距離は

$$d = \frac{|g|}{\theta} \quad (5.14)$$

である。ただし

$$\theta = \sqrt{\sum_{i=1}^n \theta_i^2} \quad (5.15)$$

これを用いて、損失を

$$l_1 = l\left(\frac{-g}{\theta}\right), \quad l_2 = l\left(\frac{g}{\theta}\right) \quad (5.16)$$

とおく。これは (4.90) で $s_\beta = 1/\theta$ とおいたものにほかならない。 $r(g/\theta)$ を計算すると、

$$r \frac{g}{\theta} = T \hat{\xi} \quad (5.17)$$

が得られる。ただし、T は行列で

$$T(\xi, \theta) = \frac{1}{\theta^3} (\theta^2 E - \check{\theta} \theta) \quad (5.18)$$

$\check{\theta}$ は θ の第 $n+1$ 成分を 0 とおいたもの

$$\check{\theta} = (\theta_1, \theta_2, \dots, \theta_n, 0)$$

である。

これより次の学習規則を得る。

$\delta\theta = \pm \varepsilon l'(|g|/\theta) C T \hat{\xi}$, ξ を誤識別したとき (+ は C_1 , - は C_2 のパターン) パターンクラスの数が多い場合にも、同様のやり方で学習規則が得られる。

識別関数を $g_\alpha(\xi, \theta) = \theta_\alpha \cdot \hat{\xi}$, $\alpha = 1, 2, \dots, k$

損失を

$$l_\alpha(\xi, \theta) = l(d_\alpha)$$

$$d_\alpha = \sum_{\beta \in N_\alpha} \frac{1}{n_\alpha} (g_\beta - g_\alpha)$$

とする*。ここに、 n_α は N_α に含まれている β の個数である。 θ は k 個のベクトル θ_α から成っているので、各成分ベクトルに分けて考えよう。すると、次の学習規則が得られる**。

一次識別関数の学習規則: $\xi \in C_\alpha$ なるパターンが誤識別されたときは

$$\begin{cases} \delta\theta_\alpha = \varepsilon l'(d_\alpha) C \hat{\xi} \\ \delta\theta_\beta = -\frac{\varepsilon}{n_\alpha} l'(d_\alpha) C \hat{\xi}, \quad \beta \in N_\alpha \\ \delta\theta_\gamma = 0, \quad \gamma \neq \alpha, \gamma \in N_\alpha \end{cases}$$

とする。

b. 線形識別関数

線形識別関数

$$g_\alpha(\xi, \theta) = \theta_\alpha \cdot \varphi(\xi)$$

の場合には、一次識別関数と同様の取り扱いができる^{42,43)}。

線形識別関数の学習規則: C_α に属するパターン ξ が誤識別されたときは

$$\begin{cases} \delta\theta_\alpha = \varepsilon l'(d_\alpha) C \varphi \\ \delta\theta_\beta = -\varepsilon \frac{1}{n_\alpha} l'(d_\alpha) C \varphi, \quad \beta \in N_\alpha \\ \delta\theta_\gamma = 0, \quad \gamma \neq \alpha, \gamma \in N_\alpha \end{cases}$$

* $s_\beta = \frac{1}{n_\alpha}$ とおいた。 $s_\beta = \begin{cases} 1, & \max_{\gamma} g_\gamma = g_\beta \text{ のとき} \\ 0, & \text{その他のとき} \end{cases}$ とおくこともできる。

** 損失を識別面からの距離の関数にとるときは、二分割の場合に行なったと同様の修正を加えればよい。

とする。

c. 区分的線形識別関数

簡単のため、二分割の場合のみを取り扱う。この場合、区分的線形識別関数は

$$g(\xi, \theta) = \max_{i=1, \dots, p} \theta_1^{(i)} \cdot \hat{\xi} + \min_{j=1, 2, \dots, q} \theta_2^{(j)} \cdot \hat{\xi} \quad (5.19)$$

と書ける。 θ は $p+q$ 個のベクトル, $\theta_1^{(i)}, \theta_2^{(j)}$ より成る。 p と q とは異なった数でよく、どちらか一方が 0 であってもよい。損失として

$$l_1(\xi, \theta) = l(-g)$$

$$l_2(\xi, \theta) = l(g)$$

をとる。

θ に対応して、学習関数 $\alpha_\alpha(\xi, \theta)$ も、各成分ベクトルに分けて考えよう。 $\theta_\alpha^{(i)}$ に対応する学習関数を $\alpha_\alpha^{(i)}(\xi, \theta)$ と書き、識別ベクトルの修正を

$$\delta \theta_\alpha^{(i)} = \varepsilon C \alpha_\alpha^{(i)} \quad (5.20)$$

とする。すると、 C_B のパターンに対する学習関数は

$$\alpha_\alpha^{(i)} = -\nabla_{\alpha^{(i)}} l_B(\xi, \theta) \quad (5.21)$$

である。ただし

$$\nabla_{\alpha^{(i)}} = \frac{\partial}{\partial \theta_\alpha^{(i)}}$$

さて、 ξ の領域 $W_1^{(i)}, W_2^{(j)}$ をそれぞれ

$$\left. \begin{aligned} W_1^{(i)} &= \{ \xi \mid \max_k \theta_1^{(k)} \cdot \hat{\xi} = \theta_1^{(i)} \cdot \hat{\xi} \} \\ W_2^{(j)} &= \{ \xi \mid \min_k \theta_2^{(k)} \cdot \hat{\xi} = \theta_2^{(j)} \cdot \hat{\xi} \} \end{aligned} \right\} \quad (5.22)$$

で定義し、

$$W_{ij} = W_1^{(i)} \cap W_2^{(j)}$$

とする。 W_{ij} においては

$$g(\xi, \theta) = (\theta_1^{(i)} + \theta_2^{(j)}) \cdot \hat{\xi}$$

である。これより $l(g)$ の微分を計算すると

$$\nabla_{\alpha^{(i)}} l(g) = \begin{cases} l'(g) \hat{\xi}, & \xi \in W_{\alpha^{(i)}} \\ 0, & \xi \notin W_{\alpha^{(i)}} \end{cases}$$

である。したがって、次の学習規則を得る。

区分的線形識別関数の学習規則：誤識別されたパターン ξ が W_{ij} にはいつているときは

$$\begin{cases} \delta \theta_1^{(i)} = \pm \varepsilon l'(|g|) C \hat{\xi} \\ \delta \theta_2^{(j)} = \pm \varepsilon l'(|g|) C \hat{\xi} \\ \delta \theta_1^{(k)} = \delta \theta_2^{(k')} = 0, \quad k \neq i, k' \neq j \end{cases}$$

(+ は C_1 , - は C_2 のパターンに対して)

とする。

多分割の場合も同様にして学習規則が得られる。

区分的学習関数を用いた簡単な実験例を示す*。パターンの特徴ベクトル ξ は二次元で、 C_1, C_2 のパターンはそれぞれ、図 5.1 (a) に示した W 字型の領域内で一様に分布している。パターンは計算機を用いて発生させる。識別関数は、4本の直線より成る区分的線形関数

$$g(\xi, \theta) = \max(\theta_1^{(1)} \cdot \hat{\xi}, \theta_1^{(2)} \cdot \hat{\xi}) + \min(\theta_2^{(1)} \cdot \hat{\xi}, \theta_2^{(2)} \cdot \hat{\xi})$$

とし、図 5.1 (b) に示される識別領域より出発して学習を行なう。識別領域の修正は、パターンを誤識別したときにのみ行なう。1回の修正で、識別領域は図 (c) のように変化し、以下 2 回めの修正で図 (d) のようになる。25 回の修正で図 (e) のようになり、ほとんど正しい識別が行なえるようになる。

* 齊藤庄司氏 (電々公社) の九州大学大学院工学研究科 (通信工学) 修士論文 (昭和 42 年) による。

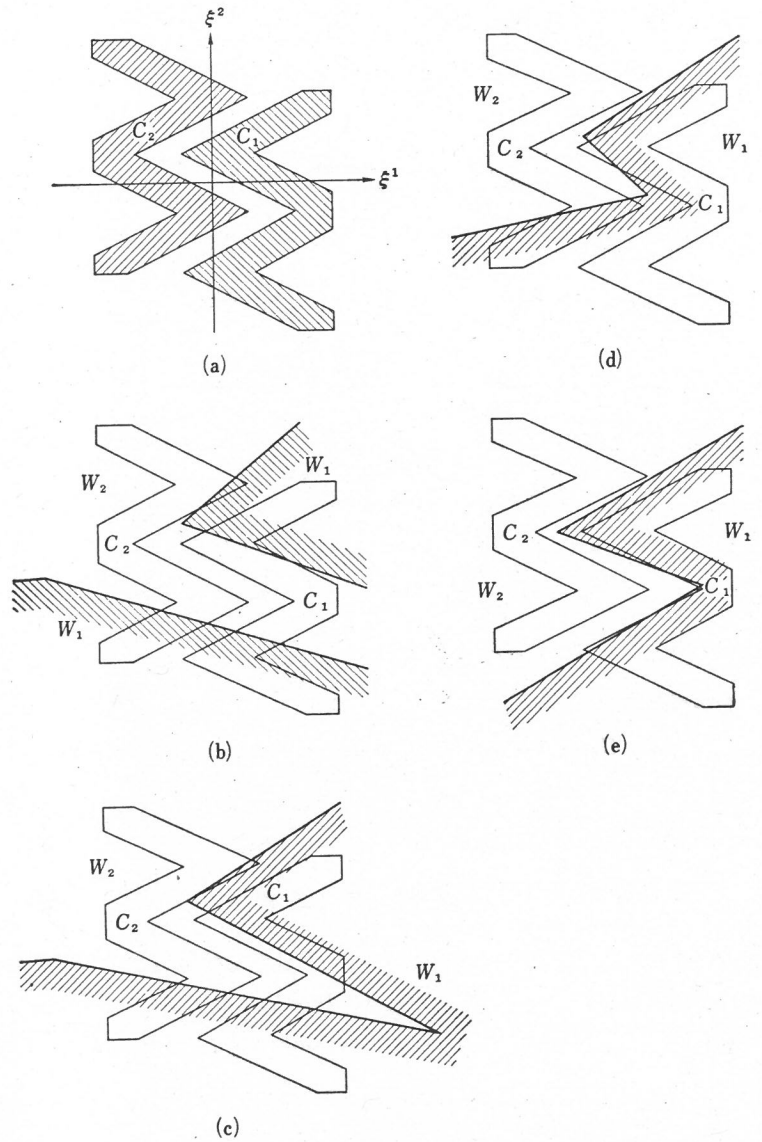


図 5.1

これは、分布が分離可能な例であるが、この学習法は分離可能でない場合でも使える。学習の収束については次節で扱う。なお、この場合でも、 $L(\theta)$ に最小値以外の極値が存在することに注意を要する。初めの識別領域の選び方によっては、正しい識別領域が得られないことがある。

この例では、四つの一次関数を用いているが、実際は三つの関数でパターンを識別することができる。事実、3回目の修正で、一つは不要になってしまう。しかし、一次関数の数を必要とするものより多めに選んでおくことは、極小値に落ち込むのを避ける点で効果がある。

C. 学習の収束

初期値 θ_1 から出発して、逐次学習を続けていく系を考える。時刻 $i(i=1, 2, \dots)$ における識別ベクトルを θ_i で表わす。これは、時刻 i に系に提示されるパターン ξ_i に基づいて修正され、

$$\theta_{i+1} = \theta_i + \delta\theta_i$$

になる。時刻 i における L の値は

$$L_i = L(\theta_i) \tag{5.23}$$

であり、これが学習によって

$$L_{i+1} = L_i + \delta L_i$$

に変わる。変化分 δL_i は時刻 i に提示されるパターンに依存している。いま、 ϵ を微小として、 ϵ^2 の項を省略すると、 δL_i の期待値は (5.6), (5.7) より

$$\overline{\delta L_i} = -\epsilon(\nabla L_i)_t C(\nabla L_i) \leq 0 \tag{5.24}$$

であることがわかる。ただし

$$\nabla L_i = \nabla L(\theta_i)$$

である。等号は

$$\nabla L_i = 0$$

すなわち、 θ_i が最適識別ベクトルであるときのみ成立する。

ここで、学習の収束を調べる。識別ベクトル θ_i は、これまでに発生したパターン $\xi_1, \xi_2, \dots, \xi_{i-1}$ の系列に依存している。パターンは確率分布 $p_\alpha, p_\alpha(\xi)$ に基づいて発生する確率変数であるから、 θ_i もまた確率変数になる。 θ_i の確率密度関数を $q_i(\theta)$ と書こう。時刻 i における L の期待値を \hat{L}_i と書くと*

$$\hat{L}_i = \int q_i(\theta) L(\theta) d\theta \quad (5.25)$$

である。1回の学習による \hat{L}_i の変化分 $\delta\hat{L}_i$ を調べると

$$\begin{aligned} \delta\hat{L}_i &= \hat{L}_{i+1} - \hat{L}_i \\ &= \sum_\alpha \int q_i(\theta) p_\alpha p_\alpha(\xi) \{L(\theta + \delta\theta(\xi, \theta)) - L(\theta)\} d\theta d\xi \\ &= \int q_i(\theta) \delta\bar{L}(\theta) d\theta \leq 0 \end{aligned}$$

すなわち、 $\delta\hat{L}_i$ は非負であることがわかる。したがって、 \hat{L}_i は単調減少する。明らかに

$$\hat{L}_i \geq 0$$

であるから、数列 \hat{L}_i は収束し、

$$\lim_{i \rightarrow \infty} \delta\hat{L}_i = 0$$

が成立する。ところが

$$\delta\hat{L}_i = \int q_i(\theta) \delta\bar{L}(\theta) d\theta$$

であり、 $\delta\bar{L}$ は最適値以外の θ に対しては常に正であるから、 $\delta\hat{L}_i$ が 0 になるためには、 $q_i(\theta)$ が θ の最適値以外のところでは、0 に収束しなければならない。これは θ_i が最適値に収束することを意味する。

以上の議論は ε^2 の項を省略したおおよっぱなものであった。収束を厳密に保証するには、 ε を定数とせず、時刻 i における値を ε_i とおいて、これをだ

* θ_i はパターン列 ξ_1, \dots, ξ_{i-1} に依存している。 \hat{L}_i は $L(\theta_i)$ をすべての可能なパターン列について平均したものである。

んだん小さくしていけばよい。この場合、修正方法は

$$\delta\theta_i = \varepsilon_i C \alpha_\alpha(\xi_i, \theta_i)$$

となる。

[定理 5.1] 条件

$$\left. \begin{aligned} \lim_{i \rightarrow \infty} \varepsilon_i &= 0 \\ \sum_{i=1}^{\infty} \varepsilon_i &= \infty \end{aligned} \right\} \quad (5.26)$$

を満たすように ε_i を選ぶならば、 θ_i は確率 1 で最適な識別ベクトルに収束する。

[証明] 確率的近似法を用いて行なえばよい⁴⁷⁾。定理の条件を満たす ε_i としては、たとえば

$$\varepsilon_i = \frac{1}{i}$$

がある。しかし、 ε_i をこのように選ぶと、収束がきわめて遅くなり、学習によって系の変動を追従する場合にはほとんど意味のないことになる。

そこで、 ε をある小さい定数に選んでおくと、どうなるかを調べていくことにする。なお、 ε の大きさを自動的に調整していく系についても、後節でふれる。

[定理 5.2] θ_{OP} を最適な識別ベクトルとし、 θ_i が θ_{OP} から μ 以上離れている確率を $M_\mu^i(\varepsilon)$ とする*。このとき、任意の μ に対して、 i が十分に大きいときは、 ε を十分に小さく選ぶことにより、 $M_\mu^i(\varepsilon)$ をいくらでも小さくできる。すなわち

$$\lim_{\varepsilon \rightarrow 0} (\lim_{i \rightarrow \infty} M_\mu^i(\varepsilon)) = 0 \quad (5.27)$$

[証明] $M_\mu^i(\varepsilon)$ を式で表わすと

$$M_\mu^i(\varepsilon) = \text{Prob} \{|\theta_i - \theta_{OP}| \geq \mu\} \quad (5.28)$$

* 最適な識別ベクトルが一つではないときは、 θ_{OP} の集合を S_{OP} とし、 $M_\mu^i(\varepsilon)$ は θ_i から S_{OP} までの距離が μ 以上離れている確率とすればよい。

となる。

$$M_\mu(\varepsilon) = \lim_{i \rightarrow \infty} M_\mu^i(\varepsilon) \quad (5.29)^*$$

とおけば,

$$\lim_{\varepsilon \rightarrow 0} M_\mu(\varepsilon) = 0$$

を証明すればよい。まず、 $\bar{\delta L}$ の厳密な式を求めておこう。 δL は

$$\delta L = -\varepsilon (\nabla L)_t C \mathbf{a}_\alpha(\xi) + \frac{\varepsilon^2}{2} \text{tr}(C \mathbf{a}_\alpha)_t \nabla \nabla L(C \mathbf{a}_\alpha) + O(\varepsilon^3)$$

と書ける。 ξ について平均をとることにすると、 $\mathbf{a}_\alpha(\xi)$ の平均は、明らかに $-\nabla L$ である。 $\mathbf{a}_\alpha(\mathbf{a}_\alpha)_t$ の平均を行列

$$Q = \int p_\alpha p_\alpha(\xi) \mathbf{a}_\alpha(\mathbf{a}_\alpha)_t d\xi \quad (5.30)$$

で表わそう。行列の関係式

$$\text{tr}(AB) = \text{tr}(BA)$$

を用いて、第2項を整理すれば

$$\bar{\delta L} = -\varepsilon (\nabla L)_t C (\nabla L) + \frac{\varepsilon^2}{2} \text{tr}(CQC_t \nabla \nabla L) + O(\varepsilon^3) \quad (5.31)$$

が得られる。したがって、十分小さい ε に対しては、定数 c が存在して

$$\bar{\delta L} \leq -\varepsilon (\nabla L)_t C \nabla L + c\varepsilon^2 \quad (5.32)$$

と書ける。

ここで

$$U_\mu = \{ \theta \mid |\theta - \theta_{OP}| \geq \mu \}$$

$$U'_\lambda = \{ \theta \mid (\nabla L)_t C \nabla L \geq \lambda \}$$

なる2種類の集合を考えよう。任意の $\mu > 0$ に対して、

$$U'_\lambda \supset U_\mu$$

* $\lim_{i \rightarrow \infty} M_\mu^i(\varepsilon)$ の収束の条件については、確率過程の本、たとえば文献⁴⁹⁾を参照。

が成り立つような $\lambda > 0$ を必ず求めることができる。

U_μ を用いて $M_\mu^i(\varepsilon)$ を表わすと

$$M_\mu^i(\varepsilon) = \int_{U_\mu} q_i(\theta) d\theta$$

となる。したがって

$$\begin{aligned} \int (\nabla L)_t C \nabla L q_i(\theta) d\theta &\geq \int_{U'_\lambda} (\nabla L)_t C \nabla L q_i(\theta) d\theta \\ &\geq \lambda \int_{U'_\lambda} q_i(\theta) d\theta \geq \lambda \int_{U_\mu} q_i(\theta) d\theta = \lambda M_\mu^i(\varepsilon) \end{aligned}$$

なる不等式が得られる。(5.32) の両辺を $q_i(\theta)$ を用いて θ について平均し、この不等式を用いると

$$\delta \hat{L}_i \leq -\varepsilon \lambda M_\mu^i(\varepsilon) + c\varepsilon^2$$

が得られる。これを i について1から N まで加え合わせ、 N で割れば

$$\frac{\hat{L}_N}{N} \leq -\varepsilon \lambda \frac{1}{N} \sum_{i=1}^N M_\mu^i(\varepsilon) + c\varepsilon^2$$

さらに、 $N \rightarrow \infty$ の極限をとると、 \hat{L}_N は有限であるから、

$$0 \leq -\varepsilon \lambda M_\mu(\varepsilon) + c\varepsilon^2$$

が得られる。したがって

$$M_\mu(\varepsilon) \leq \frac{c}{\lambda} \varepsilon$$

すなわち

$$\lim_{\varepsilon \rightarrow 0} M_\mu(\varepsilon) = 0$$

が証明される。

なお、パターン分布が分離可能であり、しかもパターンの個数が有限である場合には、有限回の学習回数で収束することがわかる。

5.3 学習の速度, 精度, 動特性

A. 学習の特性

本節では, 最適な識別ベクトル θ_{OP} の近傍における, 学習系の諸特性を調べるのに有効な補題を証明する。 θ の関数 $f(\theta)$ を考え, その値が学習の進展とともにどう変化していくかを, 一般的に考えることにする。 θ_i はパターン列 ξ_1, \dots, ξ_{i-1} に依存しているから, $f(\theta_i)$ の値もまたいかなるパターン列が過去に発生したかに関係している。時刻 i における, $f(\theta)$ の期待値を $f(\hat{\theta})_i$ で表わす。

$$f(\hat{\theta})_i = \int f(\theta) q_i(\theta) d\theta \quad (5.33)$$

以下, $q_i(\theta)$ による期待値を \hat{f}_i で表わす。

【補題】 1回の学習により, \hat{f}_i は

$$\hat{f}_{i+1} = \hat{f}_i - \varepsilon \{ (\nabla f)_t C \nabla L \}_i + \frac{\varepsilon^2}{2} \text{tr} \left(\frac{\partial^2 Q C_i \nabla \nabla f}{\partial \theta^2} \right)_i + O(\varepsilon^3) \quad (5.34)$$

に変化する。

【証明】 時刻 i における識別ベクトルを θ とし, このとき提示されるパターンを $\xi \in C_\alpha$ とする。時刻 $i+1$ には, 識別ベクトルは

$$\theta' = \theta + \delta\theta(\xi, \alpha, \theta)$$

になる。ここで θ' の確率密度関数 $q(\theta')$ を求めよう。

提示されるパターンが ξ と $\xi + d\xi$ の間にある確率は*, $p_\alpha(\xi) d\xi$ である。 $d\xi$ が

$$d\theta' = \frac{\partial \theta'(\xi, \alpha, \theta)}{\partial \xi} d\xi$$

を満足するときは, 修正された識別ベクトルは $\theta' + d\theta'$ の間にある。したがって, 修正された識別ベクトルが θ' と $\theta' + d\theta'$ との間にある確率は

* 正確には, $\xi_i \leq \eta_i \leq \xi_i + d\xi$ を満たす η_i を成分にもつパターンが発生する確率。

$$p_\alpha(\theta') d\theta' = p_\alpha(\xi) d\xi$$

である。 C_α のパターンの生ずる確率は p_α であるから, この両辺に p_α を掛けて α について加えれば, θ' の確率密度関数 $q(\theta')$ が

$$q(\theta') d\theta' = \sum_\alpha p_\alpha p_\alpha(\xi) d\xi \quad (5.35)$$

として求まる。

時刻 i における識別ベクトルを θ としたが, これは実は $q_i(\theta)$ なる分布に従う確率変数である。したがって, (5.35) の両辺の, $q_i(\theta)$ による期待値をとろう。するとこれが, 時刻 $i+1$ における識別ベクトルの確率密度を表わすことになる。したがって

$$q_{i+1}(\theta') d\theta' = \sum_\alpha d\xi \int q_i(\theta) p_\alpha p_\alpha(\xi) d\theta$$

ここで右辺の ξ は, θ と θ' との関数と考えている。

$$\hat{f}_{i+1} = \int f(\theta') q_{i+1}(\theta') d\theta'$$

であるから,

$$\begin{aligned} \hat{f}_{i+1} &= \sum_\alpha \iint f(\theta') q_i(\theta) p_\alpha p_\alpha(\xi) d\theta d\xi \\ &= \int \overline{f(\theta')} q_i(\theta) d\theta \end{aligned}$$

が得られる。ここに $\overline{\quad}$ はパターン ξ についての平均を意味する。

$$\overline{f(\theta')} = \overline{f(\theta + \delta\theta)} = f(\theta) - \varepsilon (\nabla f)_t C \nabla L + \frac{\varepsilon^2}{2} \text{tr} C Q C_i \nabla \nabla f + O(\varepsilon^3)$$

を考慮に入れれば, (5.34) が得られる。

この補題において, 関数 f はベクトル値や行列値をとるものであってもよいことに注意されたい。

B. 学習の速度と精度

前節で求めた補題を用いると, 学習の速度と精度とを求めることができる。

収束の速度は、識別ベクトルの期待値 $\hat{\theta}_i$ が θ_{OP} に近づく近づき方を示すものである。

$$f(\theta) = \theta$$

とおくと

$$\hat{f}_i = \hat{\theta}_i$$

である。したがって、補題を用いれば $\hat{\theta}_i$ に対する漸化式が得られる。

$L(\theta)$ を

$$L(\theta) = L(\theta_{OP}) + \frac{1}{2} (\theta - \theta_{OP})_t A (\theta - \theta_{OP}) + O(|\theta - \theta_{OP}|^3) \quad (5.36)$$

と展開し、 θ_{OP} の近傍のみを考えることにして、 $O(|\theta - \theta_{OP}|^3)$ の項を省略しよう。すると、学習速度の定理が得られる。

[定理 5.3] (学習速度の定理) 時刻 i における識別ベクトルの期待値は

$$\hat{\theta}_i = \theta_{OP} + (E - \varepsilon CA)^{i-1} (\theta_1 - \theta_{OP}) \quad (5.37)$$

である。

[証明]

$$f(\theta) = \theta$$

に対しては

$$\nabla f = E$$

$$\nabla \nabla f = 0$$

が成立する。これらを補題に代入すれば

$$\hat{\theta}_{i+1} = \hat{\theta}_i - \varepsilon C(\nabla L)_i$$

となる。ところが

$$\nabla L = A(\theta - \theta_{OP})$$

であるから、これを代入すると

$$\hat{\theta}_{i+1} = (E - \varepsilon CA) \hat{\theta}_i + \varepsilon CA \theta_{OP}$$

が得られる。これは $\hat{\theta}_i$ に関する差分方程式である。初期値を θ_1 としてこれを解くと、(5.37) が得られる。

これで、 $\hat{\theta}_i$ は指数的に θ_{OP} に近づくことがわかった。しかし、実際の θ_i は期待値 $\hat{\theta}_i$ に必ずしも一致するわけではない。現実の θ_i の $\hat{\theta}_i$ からのずれは、共分散行列

$$\begin{aligned} V_i &= \overline{(\theta - \hat{\theta}_i)(\theta - \hat{\theta}_i)_t}_i \\ &= \overline{(\theta \theta)_t}_i - \hat{\theta}_i \hat{\theta}_i \end{aligned} \quad (5.38)$$

で評価できる。すなわち、 V_i は時刻 i における学習の精度を表わすものとみてよい。 V_i もまた、補題に基づいて計算することができる。

[定理 5.4] (学習精度の定理) 時刻 i における識別ベクトルの共分散行列は

$$V_i = 2\varepsilon \{E - (E - \varepsilon \overline{CA})^{i-1}\} (\overline{CA})^{-1} CQC_t \quad (5.39)$$

である。特に、終局の精度は

$$\lim_{i \rightarrow \infty} V_i = 2\varepsilon (\overline{CA})^{-1} CQC_t \quad (5.40)$$

となる。ただし、 \overline{CA} は行列に対する線形演算子で、行列 M に対して

$$\overline{CAM} = CAM + (CAM)_t \quad (5.41)$$

で定義される。

[証明] ここでは、細かく式を追うことはせず、大まかな方針を示すだけに止める*。

$$f(\theta) = \theta \theta_t$$

* 証明に用いる計算に、本質的にむずかしい点はないが、相当にめんどうである。めんどうになる理由は、 f を行列とすると、 ∇f や $\nabla \nabla f$ が、3階ないし4階のテンソル量になるためである。

とにおいて、補題を用いると

$$(\widehat{\theta}_t)_{i+1} = (\widehat{\theta}_t)_i - 2\varepsilon \{CA(\theta - \theta_{OP})\theta_t\}_i^s + 2\varepsilon^2 CQC_t \quad (5.42)$$

が得られる。ここに添字 s は、行列の対称部分をとることを意味する。一方、(5.37) より

$$\hat{\theta}_{i+1} \hat{\theta}_{i+1t} = \{(E - 2\varepsilon CA)\hat{\theta}_i \hat{\theta}_{it}\}^s + 2\varepsilon (CA\theta_{OP}\hat{\theta}_{it})^s$$

が得られる。(5.42) からこれを引けば、差分方程式

$$V_{i+1} = (E - \varepsilon \widehat{CA})V_i + 2\varepsilon^2 CQC_t \quad (5.43)$$

が得られる。初期値を $V_1 = 0$ として、これを解くと (5.39) が証明される。

こうして学習の速度と精度とが求まった。いま、 CA の最小固有値を λ_0 としよう。すると、[定理 5.3] により、これに対応する固有ベクトルは収束が最も遅い方向を示し、収束の割合が時定数 $\varepsilon\lambda_0$ で示されることがわかる。一方、学習の終局の精度は $2\varepsilon(\widehat{CA})^{-1}CQC_t$ で示される。それゆえ、 ε を大きくすれば、収束は速くなるが精度は悪くなるし、 ε を小さくすれば精度は良くなるが収束は遅くなることがわかる。 A や Q についての、なんらかの知識が得られている場合には、速度と精度とのかねあいをみながら、定数 ε と行列 C を定めればよい。

C. 学習系の動特性

識別系の確率構造は時間とともに変動することが多い。この場合、最適な識別ベクトル θ_{OP} もまた変動する。学習識別系がこの変化にどう追従できるかを調べる。

時刻 i における最適ベクトルを $\theta_0 + \mathbf{A}_i$ としよう。 \mathbf{A}_i は最適識別ベクトルの変動を示している。この場合、 $\hat{\theta}_{i+1}$ の漸化式は

$$\hat{\theta}_{i+1} = (E - \varepsilon CA)\hat{\theta}_i + \varepsilon CA(\theta_0 + \mathbf{A}_i) \quad (5.44)$$

となる。行列 A も、一般には時刻とともに変動するが、ここでは簡単のため、定数行列とした。これを解くと

$$\hat{\theta}_i = \theta_0 + (E - \varepsilon CA)^i (\theta_1 - \theta_0) + \varepsilon \sum_{k=0}^{i-1} (E - \varepsilon CA)^{i-k-1} CA \mathbf{A}_k \quad (5.45)$$

が得られる。最後の項が、確率構造の変動の影響を示すものである。

最適識別ベクトルが、急に \mathbf{A} だけずれたとしよう。すると、 i 時刻後には、識別ベクトルの期待値は

$$\varepsilon \sum_{k=0}^{i-1} (E - \varepsilon CA)^{i-k-1} CA \mathbf{A} = \{E - (E - \varepsilon CA)^i\} \mathbf{A}$$

だけ変動して、このずれを追従していく。したがって

$$S_i = E - (E - \varepsilon CA)^{i-1} \quad (5.46)$$

とおけば、 S_i は学習系の階段応答 (step-response) を示すものといえる。

より直観的に応答をみるために、 θ_{OP} が周期的に変動する場合に、系がこれをどう追従するかをみよう。いま、最適識別ベクトルが $2\pi/\omega$ の周期で変動しているとして

$$\mathbf{A}_i = \mathbf{A} \sin \omega i \quad (5.47)$$

とおく。変動の周期は十分大きい、すなわち ω は微小であると仮定しておく。また、 \mathbf{A} は CA の固有ベクトルであり、その固有値を λ とする (一般の場合は、固有解を重ね合わせればよい)

$$CA\mathbf{A} = \lambda\mathbf{A}$$

(5.44) に (5.47) を代入すれば、差分方程式

$$\hat{\theta}_{i+1} = (E - \varepsilon CA)\hat{\theta}_i + \varepsilon CA(\theta_0 + \mathbf{A} \sin \omega i)$$

が得られる。この方程式の、過渡状態を示す項を除いた定常振動解は

$$\hat{\theta}_i = \theta_0 + a\mathbf{A} \sin(\omega i + \varphi)$$

とおける。これを代入すると、振幅の倍率 a と位相遅れ φ とを求める方程式

$$\begin{cases} a\{\cos(\omega + \varphi) - (1 - \varepsilon\lambda)\cos\varphi\} - \varepsilon\lambda = 0 \\ \sin(\omega + \varphi) - (1 - \varepsilon\lambda)\sin\varphi = 0 \end{cases}$$

が出てくる。

$$\alpha = \frac{\omega}{\varepsilon\lambda} \quad (5.48)$$

とおき、 α を微小とみて高次の項を省略すれば、

$$a = \frac{1}{\sqrt{1+\alpha^2}}$$

$$\tan \varphi = -\alpha$$

が求まる。それゆえ、定常解は

$$\hat{\theta}_i = \theta_0 + \frac{1}{\sqrt{1+\alpha^2}} A \sin(\omega i - \alpha) \quad (5.49)$$

である。

これは、系の周波数応答 (frequency response) を示すものである。すなわち、最適な識別ベクトルが $2\pi/\omega$ の周期で変動するとき、学習により得られる識別ベクトルは、位相が α 遅れ、振幅が $1/\sqrt{1+\alpha^2}$ 倍になって、これを追従する。したがって α が小、すなわち

$$\omega \ll \varepsilon\lambda$$

ならば、学習系はこの変動を十分によく追従できるといえる。

D. 学習法の学習

学習系の特性、特に収束速度と精度とは定数 ε と行列 C とに強く依存している。ところで、まえに示したように、 εC を調整して速度を上げるようにすると精度が落ち、逆に精度を上げれば速度が落ちる。両者の間には相反関係があって、両方を同時に良くするわけにはいかない。

しかし、よく考えてみると、速度と精度とは必ずしも同時に要求されるのではない。識別ベクトルが最適値から遠く離れているときは収束速度が速いことが望まれるが、精度すなわち平均値のまわりのばらつきはさしあたりたいした問題にはならない。ところが識別ベクトルが最適値に近づいたときは、ばらつきの小さいこと、すなわち精度が重要になってくる。したがって、識別ベクトルが最適値から遠いときは速度を速めるように εC を選び、最適値に近いとき

は精度を良くするように εC を選ぶことができれば、きわめて望ましい系が設計できる。

ところが、最適識別ベクトルは未知であるから、現在の識別ベクトルがこれに近いのか遠いのかの判断はできず、 εC をこう都合よく選ぶことはできない。そこで、 εC を過去のデータに基づいて自動的に修正していき、その結果、識別ベクトルが最適点から遠いときは速度が速まり、近いときは精度が高まるようにすることを考える。 εC は具体的な学習方法を指定する定数であるから、これは学習法の学習といえる。

識別ベクトルが最適点から離れているときは、学習により生ずる修正ベクトルは、同方向を向いているものが多く出るだろう。これに反して、最適点に近い場合には、いちど修正すると θ が最適点を飛びこすなどして、修正ベクトルの方向はまちまちになってくる。この情報を、 εC の学習に用いる。次のような εC の学習規則を考える。

識別ベクトル θ がパターン $\xi \in C_{\alpha}$ によって

$$\theta' = \theta + \delta\theta(\xi, \alpha, \theta)$$

に修正され、これがさらにパターン $\xi' \in C_{\alpha}'$ によって

$$\theta'' = \theta + \delta\theta + \delta\theta'(\xi', \alpha', \theta')$$

に修正されたとしよう。ここで

$$\delta\theta \neq 0, \delta\theta' \neq 0$$

とする (修正ベクトルが 0 である場合は、それを飛ばして、0 でないものだけを考えればよい)。このとき、次の式を用いて、 εC を $\varepsilon C + \Delta C$ に変える

$$\Delta C = \gamma a_{\alpha}(\xi, \theta) a_{\alpha'}(\xi', \theta') \quad (5.50)$$

ここに γ は正の定数である。

この修正の効果をみるために、 ΔC の期待値を求める。 ΔC をパターン ξ, ξ' について平均すると

$$\overline{\Delta C} = \sum_{\alpha, \alpha'} \int \gamma p_{\alpha} p_{\alpha'}(\xi) p_{\alpha'} p_{\alpha}(\xi') \mathbf{a}_{\alpha}(\xi, \theta) \mathbf{a}_{\alpha'}(\xi', \theta')_t d\xi d\xi'$$

が得られる。これを ξ' について積分すると、(5.9) を用いて、

$$\overline{\Delta C} = - \sum_{\alpha} \int \gamma p_{\alpha} p_{\alpha}(\xi) \mathbf{a}_{\alpha}(\xi, \theta) (\nabla L(\theta'))_t d\xi$$

となる。

$$\nabla L(\theta') = \nabla L(\theta) + \nabla \nabla L(\theta) \cdot \delta \theta(\xi, \alpha, \theta)$$

を代入して、 ξ についての積分を行なうと

$$\overline{\Delta C} = \gamma \{ \nabla L(\nabla L)_t - \epsilon Q(\theta) C_t A(\theta) \} \quad (5.51)$$

が求まる。ただし

$$Q(\theta) = \sum_{\alpha} \int p_{\alpha} p_{\alpha}(\xi) \mathbf{a}_{\alpha} \mathbf{a}_{\alpha t} d\xi$$

$$A(\theta) = \nabla \nabla L(\theta)$$

である。

識別ベクトルが最適点から遠いときは、(5.51) の第2項は第1項に比べて小さい。それゆえ

$$\overline{\Delta C} = \gamma \nabla L(\nabla L)_t \quad (5.52)$$

となる。いま ϵC を $\overline{\Delta C}$ だけ変えたとすると、修正ベクトルは

$$\overline{\Delta C} \mathbf{a}_{\alpha} = \gamma \nabla L(\nabla L)_t \mathbf{a}_{\alpha} \quad (5.53)$$

だけ変わる。

$$\overline{\nabla L \cdot \mathbf{a}_{\alpha}} = -\nabla L \cdot \nabla L \leq 0$$

であるから、修正 ΔC は、平均として $-\nabla L$ 方向を強め、収束の速度を速めるように働く。

一方、識別ベクトルが最適点に近いときは

$$\nabla L \approx 0$$

であるので、(5.51) の第1項は小さくなり

$$\overline{\Delta C} = -\gamma \epsilon Q C_t A \quad (5.54)$$

と書ける。これは、修正ベクトルを

$$\overline{\Delta C} \mathbf{a}_{\alpha} = -\gamma \epsilon Q C_t A \mathbf{a}_{\alpha} \quad (5.55)$$

だけ変える。 $Q C_t A$ は正の定符号をもつ行列であるから、この変化は、修正ベクトル $\epsilon C \mathbf{a}_{\alpha}$ を弱める方向に働く。すなわち、修正ベクトル $\delta \theta$ の絶対値が減少し、精度が上がる。

こうして、最適点から遠いときは収束速度を自動的に速め、最適点に近いときは精度を自動的に高めるような、高次の学習識別系が得られた。

一次元のパターンで、分布が重なり合っている場合について、 ϵ を自動的に変える識別実験を行なった結果では、「学習法の学習」は系の動特性をかなり高める。分布を途中で変動させた例では、人間が行なった同様の学習識別の実験結果よりも、よい結果が得られる。

この方向での学習識別理論の今後の発展が待たれる。