

直交空間での教師なし類別についての基礎研究

瀬 良 豊 士*

A Study of Unsupervised Clustering in an Orthogonal System

Toyoshi SERA

Abstract

It is known that the Karhunen-Loève system can be applied to unsupervised clustering in an orthogonal system. And some results of computer-simulation are obtained with the method for unsupervised clustering by the K-L system. But, in the case of the practical patterns, the results are not so much given. So, first the principle of the method for unsupervised clustering by the K-L system are explained here. Secondly, English character, A and B are used as the practical examples and computer-simulation by the method described above is carried out.

1. 緒 言

パターン認識の問題は、対象とするものについての予備知識があるかないかによって、次の2つに分けられる。その1つは、まえもって知られているクラス集合のいずれかに与えられた未知パターンを決定するパターン認識の問題である。他の1つは、多数の代表パターン集合から共通なもの同志を同一クラスとしてまとめクラスの集合を形成する **clustering** の問題、すなわち **cognition** の問題である。ここでは、未知パターンベクトルを直交成分に展開する **Karhunen-Loève** 法を用いて、分布の未知なパターン集合からクラスの集合を形成する **clustering** の一方法を考察する。この問題は、すでに **Watanabe**¹⁾, **Fukunaga**²⁾, **富田**³⁾ 等によって取り扱われているが、実際のパターンを用いた実験結果は余り得られていない。そこで、ここでは、まず **KL**-法を用いた教師なしの類別法が、距離の概念を用いた類別法であることを示す。つぎに、実際の未知パターン集合として、手書きの文字 **A**, **B** を用いて行なったシミュレーションの実験結果を示す。

2. Karhunen-Loève 直交系

Karhunen-Loève 展開は与えられたパターンから特徴を抽出する最も有力な方法である。ここでは、後にこの方法を教師なしの類別問題に適用するため、まず **Karhunen-Loève** 展開を導入する。この展開方法

は確率ベクトルを直交ベクトル成分の線形結合で表示する展開方法である。いま、分布が未知なる正規化されたパターン $X_{(i)}^k$ の集合 S を

$$S = \{X_{(i)}^k \mid k = 1, 2, \dots, M \quad i = 1, 2, \dots, N\} \quad \dots\dots\dots(1)$$

ただし、 $X_i^{(k)}$: m 次元パターンベクトル、

$$X_i^{(k)T} = (x_{i1}^{(k)}, \dots, x_{im}^{(k)})$$

M : パターンクラスの数。

N : 任意に抽出されたサンプルパターンの数。

とすると、パターンベクトル $X_k^{(i)}$ は

$$X_i^{(k)} = \sum_{\ell=1}^m a_{i\ell}^{(k)} \xi_{\ell} \quad \dots\dots\dots(2)$$

ただし、 $\{\xi_{\ell} \mid \ell = 1, 2, \dots, m\}$: 正規直交ベクトル。 (3)

$$(\xi^i \cdot \xi_j) = \delta_{ij}$$

$$\{a_{i\ell}^{(k)} \mid k = 1, 2, \dots, M \quad i = 1, 2, \dots, N \quad \ell = 1, 2, \dots, m\}$$

: 展開係数

$$a_{i\ell}^{(k)} = (X_i^{(k)} \cdot \xi_{\ell}) \quad (4)$$

と展開されることがすでに明らかになっている。^{1), 4)} かし、実際に $X_i^{(k)}$ を(2)式の様に展開するためには(3)の要素を求める必要がある。この要素は各クラスの生

* 工業短期大学部電気工学科

起確率を P_k とすれば, つぎの自己相関行列の固有ベクトルとして求められることが分っている.

$$G = \sum_{k=1}^M P_k \cdot E \left(X_i^{(k)} \cdot X_i^{(k)T} \right) \quad (5)$$

$$\sum_{k=1}^M P_k = 1$$

$E \left(X_i^{(k)} \cdot X_i^{(k)T} \right)$: クラス k に関する自己相関行列.

したがって, この対称行列 G の固有値 λ_l の集合を

$$\{ \lambda_l \mid l = 1, 2, \dots, m \} \quad (6)$$

とすれば, (3)の要素は

$$G \xi_l = \lambda_l \xi_l \quad (7)$$

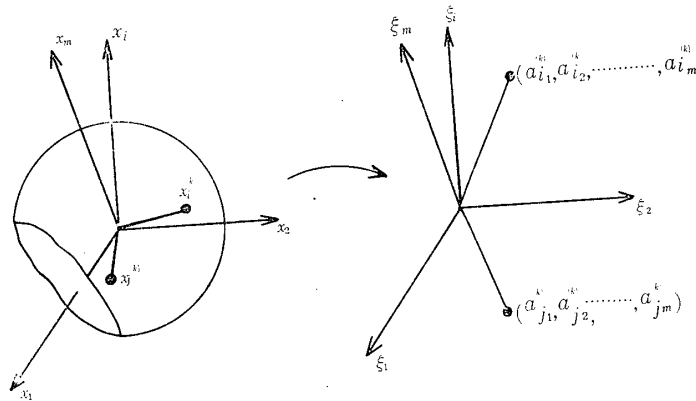


Fig. 1 Mapping two pattern vectors into the K-L orthogonal space

3. 教師なしの類別

2つの未知分布から抽出されたパターン集合を同じクラス同志に分ける2類別問題を考える. したがってパターン集合は,

$$S = \{ X_i^{(k)} \mid k = 1, 2, i = 1, 2, \dots, N \} \dots (9)$$

となる. この集合を2つのクラス Class 1, Class 2 に教師なしで類別するためには, 教師の代りとなる1つの criterion を決めておく必要がある. そこで, クラス1とクラス2の生起確率が同じであるとすれば, 各クラスの自己相関行列は

$$G_1 = \frac{1}{N} \sum_{i=1}^N X_i^{(1)} X_i^{(1)T} \quad (10)$$

$$G_2 = \frac{1}{N} \sum_{i=1}^N X_i^{(2)} X_i^{(2)T} \quad (11)$$

$$X_i^{(1)} = \sum_{l=1}^m a_{il}^{(1)} \xi_l, \quad X_i^{(2)} = \sum_{l=1}^m a_{il}^{(2)} \xi_l \quad (12)$$

となり, これらの式から G_1 行列と G_2 行列の差の行列 G_0 のトレースは,

$$t_r(G_0) = t_r(G_1 - G_2)$$

から求められる. この求められた $\{ \xi_l \}$ を Karhunen-Loève 直交系, 略して K-L という. さらに(2)式を(6)式に代入すると,

$$\lambda_l = \sum_{k=1}^m P_k \cdot Var(a_{il}^{(k)}) \quad (8)$$

$Var(a_{il}^{(k)})$: $a_{il}^{(k)}$ の分散

と求められる. この結果, (4)式と(8)式から λ_l は直交ベクトル ξ_l のパターンを表示する重要度を表わしていることが分かる. このことは $\{ \xi_l \}$ の要素を座標軸とする空間にパターンベクトルが写像された状態を考えることによってより明らかになる. この様子を Fig. 1 に示す.

$$= t_r \left(\frac{1}{N} \left(\sum_{i=1}^N X_i^{(1)} \cdot X_i^{(1)T} - \sum_{i=1}^N X_i^{(2)} \cdot X_i^{(2)T} \right) \right)$$

$$= t_r \begin{pmatrix} \frac{1}{N} \left(\sum_{i=1}^N (a_{i1}^{(1)})^2 - \sum_{i=1}^N (a_{i1}^{(2)})^2 \right) & 0 \\ 0 & \frac{1}{N} \left(\sum_{i=1}^N (a_{i2}^{(1)})^2 - \sum_{i=1}^N (a_{i2}^{(2)})^2 \right) \\ \vdots & \vdots \\ 0 & \frac{1}{N} \left(\sum_{i=1}^N (a_{im}^{(1)})^2 - \sum_{i=1}^N (a_{im}^{(2)})^2 \right) \end{pmatrix} \quad (13)$$

と求められる. これから行列 G_0 の対角要素は, $\{ \xi_l \}$ を座標とした空間を考えた場合, 各座標軸成分での2つのパターンクラス Class 1 と Class 2 の特徴差を確率平均の意味で表わしていることが分かる. このことは Fig. 1 を参照すると良く分かる. したがって, 距離の概念を導入して, つぎなる教師なしの類別の Criterion C

$$C = \sum_{l=1}^N (\lambda_l^{(1)} - \lambda_l^{(2)})^2 = \sum_{l=1}^N (\lambda_l^{(0)})^2 \quad (14)$$

$$\lambda_l^{(1)} - \lambda_l^{(2)} = \lambda_l^{(0)} : \mathbf{G}_0 \text{の固有値}$$

を採用できる. この結果 C を最大にするようにパターンを Class 1 と Class 2 に分割すれば最適な類別が行なわれることが分かる. さらに(14)式は

$$C = \sum_{l=1}^N (\lambda_l^{(0)})^2 = t_r(\mathbf{G}_0)^2 \quad (15)$$

$$t_r : (\mathbf{G}_0)^2 \text{ のトレース}$$

となる.

4. サンプルパターン交換のアルゴリズム

教師なしの類別に必要な criterion C は(14)式として求められたが, 実際にこの C を最適にもっていくためにはサンプルパターンの交換を繰り返し行なう必要がある. いま, Class 1 のパターン $X_r^{(1)}$ と Class 2 のパターン $X_i^{(2)}$ を交換したとき, 行列 \mathbf{G}_0 の変化を考えると,

$$\Delta \mathbf{G} = \frac{2}{N} (X_i^{(2)} \cdot X_i^{(2)T} - X_r^{(1)} \cdot X_r^{(1)T}) \quad (16)$$

となる. したがって交換後の \mathbf{G}_0 行列は

$$\mathbf{G} = \mathbf{G}_0 + \Delta \mathbf{G} \quad (17)$$

となる. このとき Criterion C の変化を ΔC とすると

$$\begin{aligned} \Delta C &= \sum_{l=1}^N \left((\lambda_l^{(0)} + \Delta \lambda_l^{(0)})^2 - (\lambda_l^{(0)})^2 \right) \\ &= t_r(\mathbf{G}_0 + \Delta \mathbf{G})^2 - t_r(\mathbf{G}_0)^2 \\ &= 2 t_r(\mathbf{G}_0 \cdot \Delta \mathbf{G}) + t_r(\Delta \mathbf{G})^2 \end{aligned} \quad (18)$$

となる. これから, 正しく類別するには, すなわち C を最大にもっていくには, $\Delta C > 0$ のときのみ Class 1 と Class 2 のパターンを交換することによって行なわれる. $t_r(\Delta \mathbf{G})^2$ が常に正の値であることを考えるなら

$$\Delta C = 2 t_r(\mathbf{G}_0 \cdot \Delta \mathbf{G}) \quad (19)$$

とすることが出来る. 以上より交換を行なうアルゴリズムはつぎのようになる.

- 1) $(X_r \cdot X_r^{(1)T})$ を r に固定し, i を変化させ $\Delta \mathbf{G}$ を求める.
- 2) $\Delta C > 0$ なら, この値が最大なる Class 2 のパターン $X_i^{(2)}$ を交換をする.
- 3) $r = r + 1$ として, $r = N$ になるまで(1), (2)を繰り返す.
- 4) ΔC がもはや変化しなくなったらストップする.

5) $r = N$ でも ΔC が変化するなら初めから繰り返す.

5. シミュレーション実験

集合 S のパターンとして手書文字 A, B から再構成された 9 次元の低次元パターンベクトルを採用した. 5,6) パターン数は全部で 30 個とし各クラスから任意に 15 個抽出した. したがって, シミュレーション実験が,

$$S = \{ X_i^{(k)} \mid k = 1, 2 \quad i = 1, 2, \dots, 15 \} \dots (20)$$

の条件のもとで行なわれた. 計算は ΔC が 0.005 以下になったときに停止するようにした. この結果, 9 回のパターンの交換によって類別が終了した. 類別結果は 30 個の内 2 個のパターンが誤まって類別されたのみである. Fig. 2 に criterion C の変化の様子を示す. Table 1, 2 には, 初期の状態における Class 1 と Class 2 のパターン集合, Table 3, 4 には計算終了状態における Class 1 と Class 2 のパターン集合を示す. Fig. 3 に, シミュレーション実験を行った計算機プログラムのフローチャートを示す. なお, このシミュレーション実験は TOSBAC-3400 によって行った.

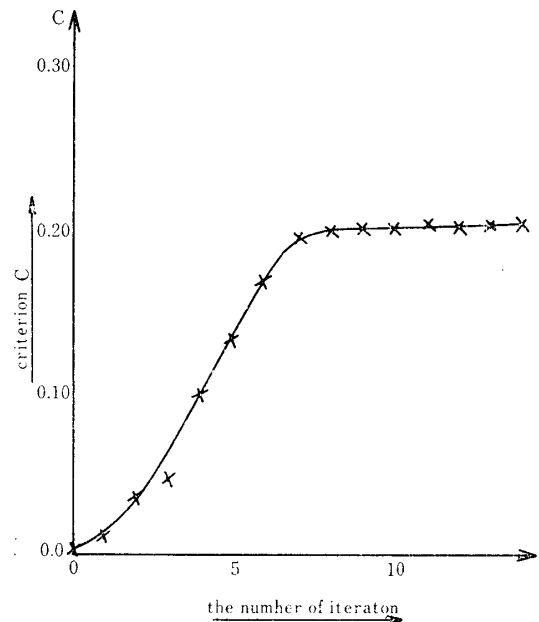


Fig. 2 Variation of criterion C against the number of iteration

Table 1 Inpnt patterns of Class 1

Dimension Known Class	1	2	3	4	5	6	7	8	9	Normalized factor
A	0.0	4.0	3.0	1.0	4.0	4.0	3.0	0.0	1.0	68.0
B	2.0	0.0	1.0	2.0	2.0	3.0	2.0	1.0	1.0	28.0
A	0.0	4.0	2.0	2.0	4.0	4.0	2.0	0.0	2.0	64.0
B	1.0	1.0	0.0	1.0	3.0	1.0	4.0	1.0	3.0	39.0
B	3.0	2.0	1.0	4.0	2.0	1.0	3.0	1.0	2.0	49.0
A	0.0	6.0	2.0	2.0	5.0	4.0	3.0	0.0	2.0	98.0
A	0.0	2.0	1.0	0.0	3.0	4.0	3.0	1.0	1.0	41.0
A	0.0	2.0	2.0	1.0	3.0	3.0	2.0	0.0	2.0	35.0
B	2.0	1.0	0.0	3.0	3.0	2.0	2.0	0.0	1.0	32.0
B	3.0	2.0	1.0	3.0	3.0	1.0	2.0	2.0	2.0	45.0
B	0.0	1.0	1.0	1.0	2.0	1.0	1.0	1.0	0.0	10.0
A	0.0	1.0	3.0	1.0	3.0	4.0	2.0	0.0	2.0	44.0
A	0.0	4.0	3.0	1.0	4.0	5.0	3.0	0.0	1.0	77.0
B	3.0	1.0	1.0	3.0	2.0	2.0	3.0	0.0	0.0	37.0
B	2.0	0.0	1.0	3.0	3.0	1.0	3.0	2.0	2.0	41.0

Table 2 Input pattens of Class 2

Dimension Known Class	1	2	3	4	5	6	7	8	9	Normalized factor
A	1.0	2.0	1.0	2.0	3.0	4.0	2.0	1.0	1.0	41.0
B	1.0	2.0	1.0	2.0	1.0	0.0	1.0	0.0	1.0	13.0
A	1.0	2.0	2.0	2.0	4.0	1.0	4.0	1.0	2.0	51.0
B	1.0	2.0	3.0	3.0	4.0	5.0	2.0	0.0	2.0	72.0
A	1.0	2.0	1.0	2.0	3.0	4.0	2.0	1.0	1.0	41.0
B	2.0	1.0	2.0	3.0	2.0	1.0	2.0	1.0	0.0	28.0
B	2.0	1.0	0.0	3.0	3.0	1.0	3.0	1.0	3.0	43.0
A	0.0	3.0	2.0	2.0	3.0	3.0	1.0	0.0	2.0	40.0
A	0.0	2.0	2.0	1.0	3.0	3.0	2.0	0.0	2.0	35.0
A	1.0	3.0	1.0	3.0	4.0	4.0	2.0	1.0	2.0	61.0
B	3.0	2.0	1.0	4.0	2.0	1.0	3.0	1.0	2.0	49.0
B	1.0	1.0	1.0	2.0	2.0	1.0	2.0	2.0	1.0	21.0
B	3.0	1.0	1.0	3.0	2.0	2.0	3.0	0.0	0.0	37.0
A	0.0	2.0	2.0	2.0	2.0	4.0	2.0	0.0	1.0	37.0
A	1.0	2.0	1.0	2.0	3.0	4.0	2.0	1.0	1.0	41.0

Table 3 The result of Calculation, Class 1

Dimension Known Class	1	2	3	4	5	6	7	8	9	Normalized factor
B	2.0	1.0	0.0	3.0	3.0	1.0	3.0	1.0	3.0	43.0
B	2.0	0.0	1.0	2.0	2.0	3.0	2.0	1.0	1.0	28.0
B	3.0	1.0	1.0	3.0	2.0	2.0	3.0	0.0	0.0	37.0
B	3.0	2.0	1.0	4.0	2.0	1.0	3.0	1.0	2.0	49.0
B	3.0	2.0	1.0	4.0	2.0	1.0	3.0	1.0	2.0	49.0
B	2.0	1.0	2.0	3.0	2.0	1.0	2.0	1.0	0.0	28.0
B	1.0	1.0	1.0	2.0	2.0	1.0	2.0	2.0	1.0	21.0
B	1.0	1.0	0.0	1.0	3.0	1.0	4.0	1.0	3.0	39.0
B	2.0	1.0	0.0	3.0	3.0	2.0	2.0	0.0	1.0	32.0
B	3.0	2.0	1.0	3.0	3.0	1.0	2.0	2.0	2.0	45.0
A	1.0	2.0	2.0	2.0	4.0	1.0	4.0	1.0	2.0	51.0
B	1.0	2.0	3.0	3.0	4.0	5.0	2.0	0.0	2.0	72.0
B	1.0	2.0	1.0	2.0	1.0	0.0	1.0	0.0	1.0	13.0
B	3.0	1.0	1.0	3.0	2.0	2.0	3.0	0.0	0.0	37.0
B	2.0	0.0	1.0	3.0	3.0	1.0	3.0	2.0	2.0	41.0

Table 4 The result of Calculation, Class 2

Dimension Known Class	1	2	3	4	5	6	7	8	9	Normalized factor
A	1.0	2.0	1.0	2.0	3.0	4.0	2.0	1.0	1.0	41.0
A	0.0	4.0	3.0	1.0	4.0	5.0	3.0	0.0	1.0	77.0
B	0.0	1.0	1.0	1.0	2.0	1.0	1.0	1.0	0.0	10.0
A	0.0	1.0	3.0	1.0	3.0	4.0	2.0	0.0	2.0	44.0
A	1.0	2.0	1.0	2.0	3.0	4.0	2.0	1.0	1.0	41.0
A	0.0	6.0	2.0	2.0	5.0	4.0	3.0	0.0	2.0	98.0
A	0.0	4.0	3.0	1.0	4.0	4.0	3.0	0.0	1.0	68.0
A	0.0	3.0	2.0	2.0	3.0	3.0	1.0	0.0	2.0	40.0
A	0.0	2.0	2.0	1.0	3.0	3.0	2.0	0.0	2.0	35.0
A	1.0	3.0	1.0	3.0	4.0	4.0	2.0	1.0	2.0	61.0
A	0.0	2.0	2.0	1.0	3.0	3.0	2.0	0.0	2.0	35.0
A	0.0	2.0	1.0	0.0	3.0	4.0	3.0	0.0	2.0	41.0
A	0.0	4.0	2.0	2.0	4.0	4.0	2.0	0.0	2.0	64.0
A	0.0	2.0	2.0	2.0	2.0	4.0	2.0	0.0	1.0	37.0
A	1.0	2.0	1.0	2.0	3.0	4.0	2.0	1.0	1.0	41.0

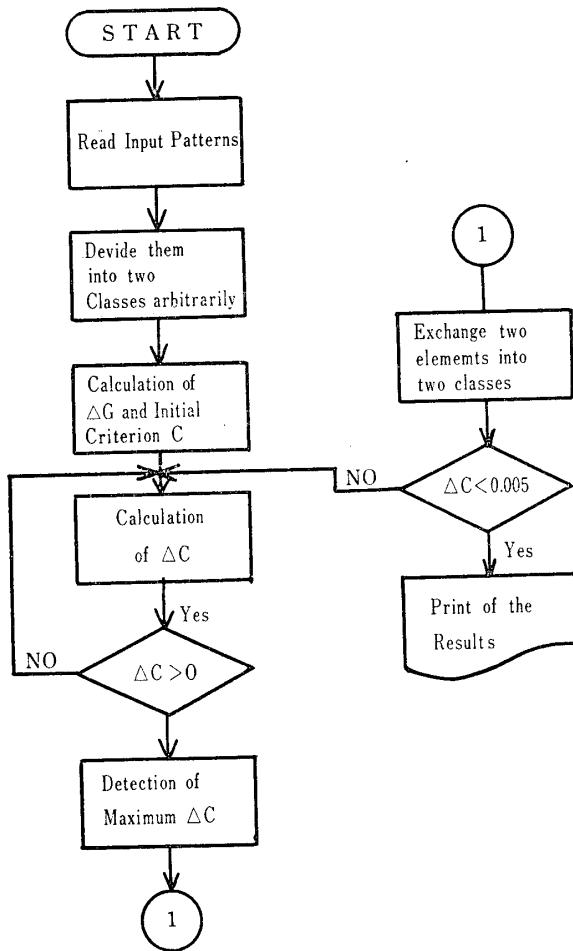


Fig. 3 Flow-Chart

6. 結 論

以上のことから

- 1) Karhunen-Loève 法はパターンの分布が未知である場合にも、教師なしで類別を行なうことができる非常に有力な方法である。
 - 2) 実験結果から(14式)の **critereion** は十分な機能をはたす。
 - 3) Karhunen-Loève 法を用いた教師なしの類別法は基本的には距離の概念を導入したものである。
- ことが明らかになった。今後は、多くの実験を行ない **C** の値が最終値に達するまでの変化の様子を検討するつもりである。最後に、日頃から何かと御世話になる工学部平田助教授に深く感謝します。同時に計算機で何かと手助けしていただいた西村女史に感謝します。

参 考 文 献

- 1) S. Watanabe: Knowing and Guessing PP 380-403 (1969) John Wiley and Sons,
- 2) K. Fukunaga, et. al: "Application of the Karhunen-Loève Expansion to Feature and Ordering" IEEE Trans. C-19 PP 314-318 (1970, April)
- 3) 富田, 他: Karhunen-Loève 直交系による教師なしの類別について, インホメーション理論研究資料 IT 70-60 (1970)
- 4) Y.T. Chien and K.S Fu: "Selection and Ordering of Feature Observation in a Pattern Recognition System" Information and Control 12 PP 395-404(1968)
- 5) 瀬良豊士: 相関法を用いたパターン認識 山口大学工学部研究報告, 23, 1, PP 1-5 (1972)
- 6) 瀬良豊士: 再構成された低次元パターンに関する考察, 電気四学会中国支部連合大会, 32a16 (1972)

(昭和48年4月14日受理)