

柔軟なオートマトンの実験的研究

東芝総研 森 健一

§ 1. はじめに

パターン認識の問題には、繁々、入力パターンが有限個の記号系列で表現され、これをいくつかのカテゴリーに識別する問題が出現する。パターンは本来多次元的な情報形態であり記号系列によつて必ずしも完全に表現されるものではないが、適当な変換によつてパターン認識の問題を記号系列の識別問題におきなおすことができる場合がある。

記号系列を識別するのに有限オートマトンを用いることが考えられるが、パターン認識の問題では入力パターンに雑音が必ず付随している入力記号系列が乱されることが多く、あるカテゴリーに属する全ての記号系列があらかじめ既知のときのみ有効な有限オートマトンの設計法ではほとんど役に立たない。また入力パターンに付随する雑音は多くの場合、入力パターンに依存関係にあるため確率オートマトンで表現す

ることもし難かしい。パターン認識における記号系列を識別するオートマトンは、Context Sensitive であり、雑音による入力記号系列の乱れに対して柔軟性をもたねばならない。また実用的な意味からは識別オートマトンの設計が容易であり、識別オートマトンの動作が速く、コンパクトに構成されねばならない。

筆者は昭和39年に自由手書郵便番号自動読取区分枝を用発する際に上記のような考察から、オートマトンに「柔軟さ」を導入することを試みた。⁽¹⁾従来の有限オートマトンでは初期状態から出発して入力記号系列によって内部状態を遷移したとき、特定の最終状態に到達したか否かによって、そのオートマトンが入力記号系列を受容するか否かを決定するが、雑音によって入力記号系列が乱されると最終状態に到達できないことが多くなる。この欠点（「固さ」）を柔らげるため、内部状態の遷移に「ペナルティ」を考え、予期しない記号が入力記号系列の途中に出現したときにはその遷移にペナルティを課し、各カテゴリー毎に作りぬに識別オートマトンの中で入力記号系列に対し最もペナルティの少なかったオートマトンがその入力記号系列を受容することにした。この場合、入力記号系列を受容したオートマトンは必ずしも特定の最終状態に到達してゐなくともよい。

各カテゴリの識別オートマトンには閾値が設けられており、もし入力記号系列が遷移の途中でこの閾値を越えるならばその時点で特定の識別オートマトンの受容が拒否されるので、各カテゴリの識別オートマトンのいずれに受容されるかを検査する手数が著しく少なくなり、速度が速くなる。

ペナルティを導入した識別オートマトンを特徴抽出過程、⁽²⁾ 識別過程に用いたところ非常に有効であることが実証された。最近に於いてこのペナルティを導入したオートマトンモデル化したペナルティ・オートマトンに関する論文⁽³⁾が発表されたのを機会に、ペナルティ・オートマトンを設計する際に必要な記号系列間の類似度と、これを用いた記号系列のラスタリングおよびペナルティ・オートマトンの設計法について考察したい。

3.2. 記号系列間の類似度

記号系列間の類似度については *Misapelin* を解析する計量言語学分野で多くの研究がなされている。⁽⁴⁾ まず予備的な定義として、記号系列とは非空の記号集合 Σ の要素を有限個並べたものを意味し、2つの記号系列 σ, ϕ があるとき

$$\begin{aligned} \sigma &= \sigma_1 \sigma_2 \cdots \sigma_i \cdots \sigma_m, & \sigma_i &\in \Sigma \\ \phi &= \phi_1 \phi_2 \cdots \phi_j \cdots \phi_n & \phi_j &\in \Sigma \end{aligned} \quad (1)$$

記号系列 $\psi = \psi_1 \psi_2 \cdots \psi_p$ が記号系列 σ および ϕ の共通部分

記号系列であるためには、少なくとも1組の添字の系列、

$\bar{i}_1, \bar{i}_2, \dots, \bar{i}_p$ および $\bar{j}_1, \bar{j}_2, \dots, \bar{j}_p$ があ、こ

$$(a) \quad p \leq \min(m, n)$$

$$(b) \quad 1 \leq k < l \leq p \text{ に対し, } \bar{i}_k < \bar{i}_l \text{ および } \bar{j}_k < \bar{j}_l$$

$$(c) \quad 1 \leq k \leq p \text{ に対し } \psi_k = \sigma_{\bar{i}_k} = \phi_{\bar{j}_k}$$

が成立たねばならない。さらに $m \times n$ 次元のマトリックス C

を考え、その要素 C_{ij} を

$$C_{ij} = \begin{cases} 1 & (\sigma_i = \phi_j) \\ 0 & (\sigma_i \neq \phi_j) \end{cases} \quad (1 \leq i \leq m, 1 \leq j \leq n) \quad (2)$$

としたとき、この C を一致マトリックスと呼ぶことにする。

一致マトリックスの要素 C_{ij} で $i=j$ の要素を a (定数) を満たす要素の集合を a 次対角と呼び、要素 C_{ij} から a 次対角への距離を $|a - (i - j)|$ で定義する。

記号系列間の類似度を定義するには、一致マトリックスの各要素に与える重み関数、一致マトリックスから最長の共通部分記号系列を選択する関数および共通部分記号系列から類似度を計算する関数を定める必要がある。⁽⁴⁾ パリ-コ認識の場合、標準記号系列に全体としてどの程度類似しているかを評価することが必要であるため、重み関数として平均的な a 次対角に対する近かきによって評価することが適切である。⁽⁵⁾

選択された共通部分記号系列の平均対角次数は $\frac{1}{p} \sum_{l=1}^p (\bar{i}_l - \bar{j}_l)$

となるので重み関数を次式で定義する。

$$W_{i_k j_k} = C_{i_k j_k} \times \left\{ 1 - \frac{1}{\max(m, n)} \left| (i_k - j_k) - \frac{1}{P} \sum_{k=1}^P (i_k - j_k) \right| \right\} \quad (3)$$

但し、 i_k, j_k は共通部分記号系列の k 番目の記号に対応するもとの2つの記号系列の記号の添字を意味する。

一致マトリックスから最長の共通部分記号系列を選択する関数は、一致マトリックスのある1の値をもつ要素 C_{ij} に対応する記号を選んだとき、次の記号として $k > i, l > j$ で $C_{kl} = 1$ となる要素に対応する記号を選んで構成した共通部分記号系列の中で最も長い系列を選ぶ。

例1. $\Sigma = \{A, B, C, D, E, F\}$

$\sigma = ABCDEF$

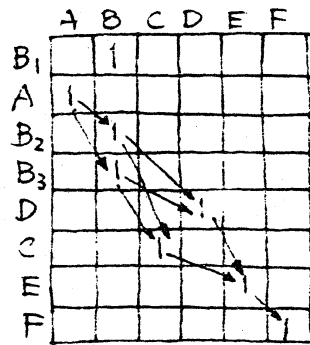
$\phi = BABBDCEF$

とすると最長の共通部分記号系列は4通りある。それを E であらわすと

$E = \{AB_2DEF, AB_3DEF, AB_2CEF, AB_3CEF\}$.

ある最長の共通部分記号系列 ψ に対し、その構成要素の一致マトリックスよでの重み関数を (3) 式とすると、2つの記号系列 σ および ϕ に関して共通部分記号系列 ψ を介しての類似度を次式で定義する。

$$S(\sigma, \phi / \psi) = \left\{ 1 - \frac{\frac{1}{2}(m-n) - \frac{1}{P} \sum_{k=1}^P (i_k - j_k)}{\max(m, n)} \right\} \cdot \sum_{k=1}^P \frac{W_{i_k j_k}}{\max(m, n)} \quad (4)$$



一致マトリックス

さらに2つの記号系列 σ および ϕ の類似度を次式で定義する。

$$S(\sigma, \phi) = \max_{\psi \in E} \{ S(\sigma, \phi / \psi) \} \quad (5)$$

2つの記号系列間の類似度 $S(\sigma, \phi)$ は

$$0 \leq S(\sigma, \phi) \leq 1 \quad (6)$$

とほり、記号系列 σ と ϕ の共通部分記号系列が空のとき0となり、記号系列 σ と ϕ が一致するとき1になる。

例2. 例1の問題に対し、類似度を計算する。

$$\psi_1 = AB_2DEF, \psi_2 = AB_3DEF, \psi_3 = AB_2CEF, \psi_4 = AB_3CEF$$

$$S(\sigma, \phi / \psi_1) = 0.558$$

$$S(\sigma, \phi / \psi_2) = 0.543$$

$$S(\sigma, \phi / \psi_3) = 0.517$$

$$S(\sigma, \phi / \psi_4) = 0.519$$

$$\therefore S(\sigma, \phi) = 0.558$$

上の定義では記号間の類似度は1か0で評価した。すなわち、 $A, B \in \Sigma$ とすると $S(A, A) = 1$. $S(A, B) = 0$. しかしながら、パターン認識の問題では記号間の類似度が全て1か0になる場合は少なく、異なる記号間の類似度が適当な値をもつ場合が多い。この場合に記号系列間の類似度の定義を拡張することは容易で、まず一致マトリックス

の要素 C_{ij} を記号系列 σ の i 番目の記号 σ_i と、記号系列中の j 番目の記号 ϕ_j の類似度で置直し、(2) に代る

$$C_{ij} = S(\sigma_i, \phi_j), \quad (1 \leq i \leq m, 1 \leq j \leq n) \quad (2)'$$

とする。共通部分記号系列の選択は 0 でない要素の組合せの全てについて行なう必要がある。重み関数と類似度の評価式は (3), (4), (5) 式を用いる。

例 3. 例 1 に $S(A, B) = 0.6$, $S(C, D) = 0.7$ を加えた場合の記号系列 ABCDEF と BABBDCEF の類似度を計算する。

$$S(\sigma, \phi / AB_2DCEF) = 0.583$$

$$S(\sigma, \phi / AB_2DEF) = 0.558$$

$$S(\sigma, \phi / B_2B_3DCEF) = 0.523$$

$$S(\sigma, \phi) = 0.583$$

$\sigma \backslash \phi$	A	B	C	D	E	F
B ₁	0.6	1				
A	1	0.6				
B ₂	0.6	1				
B ₃	0.6	1				
D			0.7	1		
C			1	0.7		
E					1	
F						1

§3. 記号系列のクラスタリング

識別オートマトンを設計する際にあるカテゴリーに属する記号系列のサンプルが与えられたとき、サンプル全体を代表する記号系列はどれかをきめたいときがある。すなわちあるカテゴリーに属する標準記号系列をサンプル記号系列群からきめようとする問題がある。さらにサンプル全体を1つの記

号系列で代表させることが困難なときは、サンプル記号系列群を適当な数の類似群にクラスタリングする問題を考えることができる。

今サンプル記号系列の集合を G とし、 G に属する記号系列の数を N とするとき

$$S_{\sigma} = \frac{1}{N-1} \sum_{\substack{\phi \in G \\ \phi \neq \sigma}} S(\sigma, \phi), \quad \sigma \in G \quad (7)$$

を計算し、 S_{σ} が最大となる記号系列をサンプル記号系列 G の代表記号系列といる。

サンプル記号系列 G を任意の数 M 個の類似群にクラスタリングする手順は、多次元ベクトルで表現されたサンプル系列群をユークリッド距離を評価関数として漸近的にクラスタリングする手法⁽⁶⁾に準じて、次の通りとすることができる。

記号系列群のクラスタリングの手順

- (1) サンプル記号系列の集合 G から適当な M 個の記号系列を核記号系列として選ぶ、
- (2) G に属する任意の記号系列をとり出し、核記号系列との類似度を計算し、最も大きな類似度を示す核記号系列の群にその記号系列は属するとして分類する。 G の全ての記号系列を分類し、 M 個の群に分ける。
- (3) 各分類群の代表記号系列を計算し、これを新しい核記

号系列とする。

- (4) 各分類群の代表記号系列が変化しなくなるまで(2)(3)の手順を繰返す。もしある分類群で複数回の記号系列が同じ最大の S_0 をもつときは、他の分類群の代表記号系列に平均的に最も類似していい記号系列を代表記号系列とする。

例4 $G = \{BA, BB, BBA, ABB, ABA, AABA\}$ の代表記号系列を求め、次に2つの類似群に分類する。

まず記号系列間の類似度を計算すると右表のまうにたり、これから代表記号系列は ABA である。最も類似していい BB と $AABA$

	BA	BB	BBA	ABB	ABA	AABA
BA	1.0	0.5	0.56	0.28	0.50	0.38
BB	0.5	1.0	0.56	0.56	0.28	0.25
BBA	0.50	0.56	1.0	0.44	0.67	0.44
ABB	0.28	0.56	0.44	1.0	0.67	0.44
ABA	0.50	0.28	0.67	0.67	1.0	0.60
AABA	0.38	0.25	0.44	0.44	0.60	1.0

を最初の核記号系列とすると $(BB), (ABA, BA, BBA, ABB, AABA)$ の2群に分類される。最も類似した ABB と ABA を最初にとっても同じ結果が得られる。しかし、 BBA と ABA を最初にする $(BBA, BA, BB), (ABA, BA, ABB, AABA)$ に分割される。

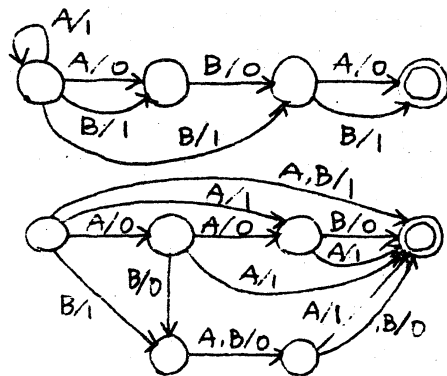
最初に与える核記号系列によっても常に同一のクラスタリニクがなされるかが、分類された結果はほとんど似た類似群が構成されるという意味で安定したクラスタリニクである。

§4. パナリティ・オートマトンの設計法

文献(3)のヤナルティオートマトンの構成法をヒントにし、記号系列間の類似度を用いて並列ヤナルティオートマトンを設計する手順は、次のようになる。

- (1) 各識別カテゴリー毎にサンプル記号系列から代表記号系列を求め、ヤナルティ、0のヤナルティ、オートマトンを構成する。
- (2) 代表記号系列との間の類似度が大きい順にサンプル記号系列を受容するようにヤナルティ、1の遷移を付加していく。同時に他のカテゴリーの記号系列を受容しないようにチェックする。
- (3) あるカテゴリーに属する記号系列が自分の代表記号系列より他のカテゴリーの代表記号系列により類似しているときは、そのカテゴリーを2つにクラスタリングし、新しい遷移ルートをヤナルティ、1で付加する。

例5 (BA, BB, BB \bar{A} , \bar{A} BB, \bar{A} BA, \bar{A} BA \bar{A}), (A, B, AA, \bar{A} B, AAA, AAB, BAA, BAB, BBB, \bar{A} BB \bar{A}) をサンプル記号系列とする並列ヤナルティオートマトンを構成する。



§5 おわりに

記号系列の識別オートマトンを構成するのに必要な記号系列間の類似度および記号系列のリスタリングの手順について述べ、並列パナルテ、オートマトンの設計への応用を試みる。並列パナルテ、オートマトンの理論や構成法についてはまだ多くの研究すべき事項が残っている。

参考文献

- (1) H. Genchi, K. Mori, S. Watanabe and S. Katsuragi, Proc. IEEE, 56, 8, pp 1291-1301, (1968).
- (2) K. Mori, H. Genchi, S. Watanabe, S. Katsuragi, Pattern Recognition, 2, pp 175-185, (1970).
- (3) 阿部圭一, 電通学会研究会資料. PRL74-5, pp 45-54, (1974).
- (4) C. N. Alberg, U. ACM, 10, 5, pp 302-313, (1967)
- (5) R. Vivès and J. Y. Gresser, Proc. IJCPPR, pp 308-317, (1973).
- (6) G. H. Ball, Proc. FJCC, Vol 27, (Spartan Books), pp 533-559, (1965).