

系列パターン認識システムの考え方 — 系列パターンの誤り処理 —

京大 工学部 堂下 修司
山崎 進

引 まえがき

パターン認識の問題において、刺激として与えられる位置パターンと、最終的な識別結果の論理的模式との間、通信チャンネルとして系列パターンが存在する。位置パターンや論理的模式については、信号として、または論理的関係としての近さが存在し、それに基づいて、パターン認識が行なわれるのであるが、系列パターンに関しては、そのような近さの尺度を導入することは、必ずしも自明ではない。一般的には、この系列パターンは言語の問題として理論的に扱うことが可能である。系列パターンの分類を言語の問題として扱う場合、文法あるいはオートマトンを用いて、系列パターンの各クラスに対して文法を構成し、その各々のクラスの識別を、文法による系列の受理非受理で行うこととなるか

あるクラスに対して、それに属する系列を受理するとしても、そのクラスの系列パターンのもつてゐる構造上の関係を正しく反映する文法あるいはオートマトンを構成すること自体の問題であり、オートマトンの構成に関し⁽¹⁾ 複本代らのグループによつて詳細な研究が行なわれてゐる。

一方、ある文法と言語が与えられてゐるときに、その言語に属する記号列に対して、実際に観測される系列パターン入りに誤りを含む場合、それを与へた言語から生成されたものとして認識する(構文を求める)問題は類似の問題ではあるが上記の系列パターンの分類の問題とは異なつた観点からの取扱ひが必要となる。すなわち、ある言語を正しく受理する文法が与えられてゐる場合、その言語に属する記号列が誤りを含んだパターンとして得られたとき、文法から、その誤りを検出し訂正する方法が実際的課題となる。とくに、系列を構成する記号が独立に、記号間の置換、記号の挿入、記号の削除によつて誤るといふ条件が最も簡単な場合として考えられるが、正規言語より上のクラスの言語に対しては、有限個の誤りに関して、その誤り検出が可能かどうかを、一般的には決定できず、誤り系列の訂正に関してもあつまい性が残る。⁽²⁾

したがつて、ある種の制限下での誤り訂正、例えば最小誤

りの誤り訂正を行うというところが現実的課題となり、そのためのアルゴリズム手法が提案されている⁽³⁾。

パターン認識の立場にかゝるとは、系列パターンが、ある確率をもつ長誤りの構造の下で生起してゐると見るのが自然である。とくに記号の置換によつて誤りが確率的に生じてゐるとするときは、入力記号列に対し確率的に見て最も「確からしい」ものへの訂正方法が述べられている⁽⁴⁾。

その他、文法による系列パターンの分類にかゝると、系列パターンの生成および誤り変形を確率的に見て、系列パターンの識別に確率的文法を用い *stochastic syntax analysis* を適用し、染色体の識別に応用した例が論じられている⁽⁵⁾。

我々は、系列パターンを言語論的に扱う上で、言語に属するそれそれの系列パターンの生起確率が予じめ提示され、入力系列パターンの誤り生成が確率的にわかつてゐるとし、系列パターン集合の系列の生起確率から見た特性を反映した文法を構成することおよび、入力記号列パターンのもとでの言語への同定という問題に興味をもつてゐる。ここでは、記号列集合 L と、それを特性づける確率的文脈自由形文法 G が、何らかの方法で与えられてゐるとし、実際に観測される誤りを含んだ系列パターン集合 $L'(C, L)$ とするときは、入力記号列パターン $w' (\in L')$ の記号列 $\{w\} \subset L$ への同定を、記号列 w の生

起確率および w が w' に誤る確率の積が最も大きいような w を求めるいわゆる最尤推定の問題として考え、 $\{w\}$ を求める手法としては、文献(3)の手法を変形すれば、効率的なものが得られることを示す。

§2. 系列パターン同定の

まず若干の準備を行う。文脈自由形文法(cfg) $G=(V_N, V_T, P, S)$ 等は通常の記法に従う。確率的自由形文法(scfg)⁽⁶⁾ $G_s=(G, \varphi)$ である。 G は cfg であり、 φ は P から実数区間 $(0, 1]$ への写像、 $\varphi: P \rightarrow (0, 1]$ である。ただし $\varphi(P_i), P_i \in P$ を同一の左辺に $\varphi(P_i)$ を加えて $\varphi(P_i)$ の和が 1 に等し... とする。 $w \in L(G)$ に対し、 $\psi(w) = \sum_{P_1 P_2 \dots P_n \in E(w)} \varphi(P_1) \varphi(P_2) \varphi(P_3) \dots \varphi(P_n)$ ($E(w)$ は w の左生成過程の全体) によって w の生成確率を与える。 $L(G_s) = \{(w, \psi(w)) \mid w \in L(G)\}$ とする。

ここでは簡単の場合としてつぎのような誤りを考える。系列に対して、系列を構成する記号の誤りが独立に、記号間の置換、記号の削除、記号の挿入によって生ずるものとし、 $a \in V_T^U \setminus \{ \epsilon \}$ が $b \in V_T^U \setminus \{ \epsilon \}$ に誤る確率を $r(b|a)$ と書く。 $a \in V_T$ が誤らないときも $r(a|a)$ なる確率での一種の誤りと見なす。 $r(b|a)$ が図1のような形で予め定義されるものとする。($r(\epsilon|\epsilon)$ も形式的に定義される。) 系列 $x = a_1 a_2 \dots a_n$ が $a_1 \rightarrow b_1, a_2 \rightarrow b_2, \dots, a_n \rightarrow b_n$ と誤ることによって系列 $y = b_1 b_2 \dots b_n$ が得られると

したがって誤り確率 $g(y|x)$ は $g(y|x) = r(b_1|a_1)r(b_2|a_2)\dots r(b_n|a_n)$ 与えられるとする。

誤り 記号	V_T	ϵ
V_T	誤り 置換	置換 削除
ϵ	削除	誤り 置換

$a \in V_T$ に対し

$$\sum_{b \in V_T \cup \{\epsilon\}} r(b|a) = 1$$

図1 入力系列ホターンに関する誤りの構造

いま、系列ホターンの誤り処理を次のような問題と考える。「 $G_S = (G, \varphi)$, $G = (V_T, V_n, P, S)$ が定義されてゐるとき、 $w \in V_T^*$ に対し $\max_{w' \in L(G_S)} \varphi(w') g(w'|w)$ なる w' の集合を求め、 w' の構成を得る手順を考えること」

以下、手順としては、文献(3)の手

法を変形すれば通用できることを示す。

文献(3)では Earley の アルゴリズムを用いる、解析中のフォークシヨの番号とその右辺の参照位置および入力系列上の参照位置の三つ組から成つてゐる“状態”に誤りの回数を表わす変数を組み込んでゐるが、その代りに $(0,1]$ の確率を組み合わせることとする。概略を示すと次のようになる：

1) G_S に $Z \rightarrow S^+$ を付け加える。 G_S の生成規則に番号付けを行ひ、 $Z \rightarrow S^+$ のそれを 1 とする。入力は $w^+ = t_1 t_2 \dots t_m^+ \in V_T^+$ とする。($t_i \in V_T$)

2) (P, j, f, e) [P : 生成規則の番号, j : P の生成規則の右辺の位置を示す, f : $t_1 t_2 \dots t_m^+$ のある位置を示すバックポイント

$\gamma, e: \text{実数値}$] を状態とする。

3) 状態の集合 $\text{statesets } \{S(i)\}$ を考え、入力の位置 $t(i)$ に対応して $S(i)$ とする。 $S(i)$ に後述するように 4 つ組の状態を付け加え之の中。 $w^{-1} = t(1) t(2) \dots t(m)$ に対し $m+2$ の statesets がある。 Statesets へ 状態 (p, j, f, e) を加えるとき、その中 (p, j, f, e') ($e > e'$) があれば (p, f, j, e') は (p, f, j, e) におきかえられる。

4) 処理の手順 ① $S(i)$ は $(k, 1, i, 1)$ $1 \leq k \leq \#P$ と初期設定される。 ② $S(i)$ から $S(i)$ の state を k だけおきかえられる。 ③ $S(i)$ の後 Completer は $S(i)$ の全ての finalstate をおきかえる。 \therefore finalstate とは状態の 4 つ組において k 番目の生成規則 P_k に対し、 j が P_k の右辺の最後を示しているときの状態である。 ④ $S(i)$ から ②, ③ の処理が行われれば、 $S(i+1)$ への処理 ②, ③ を繰返し、 $S(m+2)$ で終了する。

S(i) ① \therefore \therefore での入力 $t(i)$ の check と推定処理を行う。 ② 入力 $t(i)$ に対し \therefore ぎのようになる。 (i) $t(i)$ が P 番目の生成規則の左辺の j 番目 $C(P, j)$ に対し (p, j, f, e) ($j+1 \in P(i)$) で一致すれば $(p, j+1, f, e \cdot r(t(i) | t(i))) \in S(i+1)$ 。 $(p, 1, f, e) \in S(i)$ に対し $t(i)$ が $C(P, 1)$ と一致すれば $(p, 2, f, e \cdot r(t(i) | t(i)) \cdot \varphi(P)) \in S(i+1)$ ($\varphi(P)$ は P の適用確率) (ii) $(p, j, f, e) \in S(i)$ ($j+1 \notin P(i)$) に対し $C(P, j) \neq t(i)$, $C(P, j) \in V_T$ ならば $(p, j+1, f, e \cdot r(t(i) | C(P, j))) \in S(i+1)$ 。 $(p, 1, f, e) \in S(i)$ に対し $C(P, 1) \neq t(i)$, $C(P, 1) \in V_T$ ならば $(p, 2, f, e \cdot r(t(i) | C(P, 1)) \cdot \varphi(P))$ 。 (iii) $(p, j, f, e) \in S(i)$ ($j+1$)

で、 $C(P, j) \in V_T$ ならば $(P, j+1, f, e \cdot r(\epsilon | C(P, j))) \in S(i)$. $(P, 1, f, e) \in S(i)$ で $C(P, 1) \in V_T$ ならば $(P, 2, f, e \cdot r(\epsilon | C(P, 1))) \in S(i)$. (iv) $(P, j, f, e) \in S(i)$ ならば $(P, j+1, f, e \cdot r(t | \omega | \epsilon)) \in S(i)$ とする。

Compiler ① $S(i)$ の上で SCAN が終了したとき、働く処理ルーチンである。② SCAN で $C(P, j)$ が V_N かつ r があるときは、 $C(P, j)$ を左辺とする生成規則をもつ状態を $S(i)$ に加える。③ final state を次の例のような手順をおきかえる。 ($V_p \rightarrow W \cdot V, e$) $\in S(i)$ (j は a の位置によって表わされる。) ($V_x \rightarrow r, e'$) $\in S(i)$ のとき ($V_p \rightarrow W \cdot V \cdot a, ee'$) $\in S(i)$ とする。可能な部分系列の構文の完了を意味する。この構文の部分木は適当に記憶しておく。

最後に $(Z_1 \rightarrow S^{-1}, e) \in S(n+2)$ が得られたとき、 $e = \max_{w \in L(G_S)} P(w|w)$ であり、 w の構文木を得られる。

次に例を示す。

$$G_S = \{ \{s\}, \{a, b\}, \{s \rightarrow asb, s \rightarrow b\}, s, \{ \varphi(s \rightarrow asb) = 0.6, \varphi(s \rightarrow b) = 0.4 \} \}$$

とし、誤りの構造を図2のように仮定する。

$w \backslash$	a	b	ϵ
a	0.7	0.2	0.1
b	0.1	0.7	0.2
ϵ	0.1	0.1	0.8

一般的に r 対角線要素の値は、他の要素の値に比べて大きく、また一定以下の確率

図2.

での誤り(例えば0.1以下)はゼロとみなして無視することが出来る。これにより、入力系列パターンへの信頼性が高く、あさましさが少い場合、構文解析への負担が軽くなり、平等な全々の誤りを可能とみて e を誤り個数として計算する場

金に比べて利点があると思われる。このとき, Statesets は図3
 のようになる。図において, M, S, D, I はそれぞれ SCAN によつ
 て, マッチング のとれ長ニヒ, 置換, 削除, 挿入 による誤り
 があるとして処理するニヒに対応している。4つ組の状態の
 第1番目の生成規則において, $(P, j) \in V_N$ から出てくるポイント
 は Completer による。①から④までは, 最終状態のおきかえで

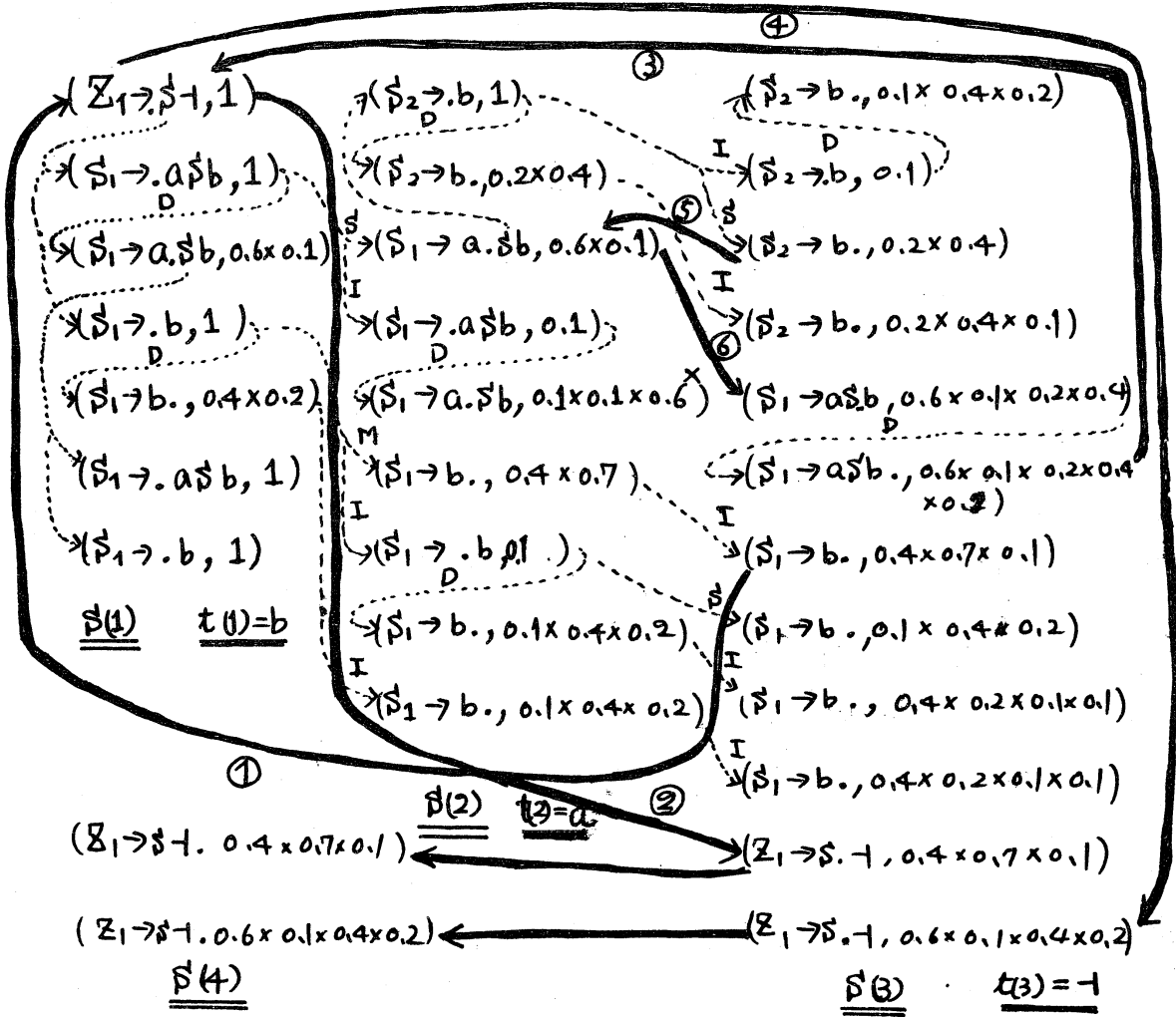


図3. 文法 G_3 , $\lambda b:ba$

ある。この単純な例では、 $b \in L$ が $\max_{w \in L(G)} \psi(w) \delta(w|w)$ を満たす w に付しておりの挿入誤りが生じたと判断している。上記の手順においては、 $P(w, w) = \max_{w \in L(G)} \psi(w) \delta(w|w)$ を満たす w のみを考えたい。しかしこの最良の判定における $P(w, w)$ の値が識別結果として十分大きく、 $P(w, w)$ が他の $P(w, \hat{w})$ ($\hat{w} \in L, \hat{w} \neq w$) に対する値に比べて十分大きい場合には、これが良いが、 $P(w, w)$ の値が不満足であつたり、 $P(w, w)$ と $P(w, \hat{w})$ の比が接近している場合、一意的に w を w と見做すことは問題がある。実際には、 $\psi(w) \delta(w|w)$ の値が大きい方からいくつかの記号系列 w を候補として列挙することが望ましい。特にこのような判定を用いて、さらに上位の識別の過程（たとえばセマンティクス）などを考える場合、最も確からしいいくつかの候補をその判定の値 P と共に求めておくことは必要であると思われる。

§3 むすび

系列パターンを文法によつて扱うとき、系列パターンが記号の間にある確率をもつて独立に生ずる誤りを含んで生成される場合に、系列パターンを最も高い確率で同定するのに必要の手順は、言語の最小誤り訂正の手法において用いられるエラーカウンターを別の量で定義すれば、その手法と同等の記憶量、解析時間が可能であることを検討した。

今後の問題として、ある言語に属する記号列の生起頻度が

予じめ知られている場合, それを規則の確率または, 規則の列に対する条件は, 確率として近似的に文法に賦与する方法も考える必要がある。

文献

- 1) 榎本, 堂下, 富田: カニフォル記号列を識別する最簡オートマトンの構成, 信学論 vol 55-D, No.3, PP 210-217, 1973年3月
- 2) 岩元, 沢野: 正則言語, 文脈自由言語の誤り訂正, 信学論 vol 56-D, No 12, PP 675-680, 1973年12月
- 3) G. Lyon: Syntax-directed least errors analysis for Context-free languages: A practical approach, CACM. vol 17, No2, PP3-14 Jan.'74
- 4) L. W. Fung and K.S. Fu: Stochastic syntactic classification of noisy patterns, Proceedings of the 2nd International Joint Conference on Pattern Recognition, pp102-103, '74
- 5) H.C. Lee and K.S. Fu: A stochastic syntax analysis procedure and its application to pattern classification, IEEE Trans. on Computers, vol C-21, No.7, PP 660-662, July, '72
- △) 尾関: 文脈自由形言語の情報論的性質, 信学論, vol 55-D, No.11, PP 753-760, 1972年11月
- 7) T. L. Booth: Probabilistic representation of formal languages, IEEE Trans Annual Symposium on Switching and Automata Theory, PP 74-81, Nov. '69