

情報圧縮と文法推論

九大 理 情報研 有川節夫

文法推論による情報圧縮とその可能性について考へ、
さらに、右延長文法とその推論について考へらる。

1. 情報の圧縮と損失

一般に情報処理の過程はある目標 θ に沿った何段階かのデータの変換からなるといえる。データ D を記録するのに要する面積を $|D|$ とするとき、 $|D| > |T(\theta, D)|$ であれば、 D は変換 T によって圧縮されたことになる。変換には、検索効率を高めるためにデータを組織化する場合、あるいは、データを図式化する場合のように $|D| \leq |T(\theta, D)|$ となる変換もあるが、多くの変換は圧縮と指向したものであるといえよう。

また、変換には、とくに、圧縮の場合には、“情報の損失”と伴うことが多い。ここで情報の損失の定義が問題になるがそれは、情報処理過程に対する目標 θ を介してつぎのように捉えられる：即ち、変換 T が目標 θ の下でデータ D に対して情報の損失を伴わないとは、 (θ, D) に適用可能な任意の変

換 R に対して

$$R(\theta, D) = S(T(\theta, D)) \quad (1.1)$$

となる変換 S が存在することである。こうすると、一見大幅な圧縮をなす変換でもその目標のためには、情報の損失は無いことがあることなども説明できる。

一般には、情報の損失を伴わない大幅な圧縮は目標が単純明確である場合には可能であることが多いが、多目的に使用されるファイルの作成、あるいは用途の細目が完全には分っていないデータの整理などの場合には困難である。

2. 文法推論による圧縮

このような局面に対処する1つの方法として、 θ に依存しない情報の圧縮、損失と考へる、即ち、“任意の θ の下で情報の損失を伴わない圧縮”を考へる方法がある。この場合には、変換は単に $T(D)$ の形で扱われ、情報の損失に属する式 (1.1) は

$$D = S(T(D)) \quad (2.1)$$

で置き換えられる。“ T が D に対して情報の損失を伴わない”とは“ $T(D)$ から D を復元することが可能である”ということになる。このように θ を無視することにより、圧縮率の減少はみぬかれないうえ、一般的な議論の可能性が与えられる。

一般的な1方法として、ここでは、文法推論の方法を利用

あることを考へる。文法推論とは、与えられたサンプル集合から、その特徴を抽出して、それを文法あるいはオートマトンとして記述する方法である。

ところで、データは記号列の系列または集合として扱われるが、ここでは単に記号列の有限集合であると考えらる。即ち、 Σ を記号の有限集合とすると、データ D は Σ^* の有限部分集合であるとある。

文法推論による圧縮法とは、データ D に対して、文法族 $G = \{G_1, \dots, G_n\}$ と述語 P とを条件

$$D = L(G, P) \quad (2.2)$$

$$|D| > |G| + |P| \quad (2.3)$$

を満たすように決定する方法である。ここで、文法族、述語の選択範囲、 L の定義が問題になる。文法に関しては、効率より推論が行えるために、有限オートマトンかそれに類したもの、述語はアトリビュートの *rudimentary attribute* 的なものが妥当であろう。 L はその際に自然に定まる。

しかし、当面は現在の文法推論法を利用するために、

$$D \subseteq L(G(D)) \quad (2.4)$$

なる文法 $G(D)$ を推論して、適当なフィルターとしての述語 P を選り、余分な $L(G(D)) - D$ と振り落して

$$D = L(G(D), P) \equiv \{x \in L(G(D)) ; P(x)\} \quad (2.5)$$

とあることを考えよう。たとえは、データ

$$D = \{a, aab, aabab, \dots, a(ab)^{100}\} \quad (2.6)$$

に対して、文法と述語

$$G(D) = \{S \rightarrow aA, A \rightarrow abA, A \rightarrow \varepsilon\} \quad (2.7)$$

$$P(x) \leftrightarrow |x| < 202 \quad (2.8)$$

を指定すれば、 $D = L(G(D), P)$ となる。この場合の変換 T は、 D から $G(D)$ と P を作る操作である。データの面積は区切り記号を無視すれば、 $|D| = 10201 > 19 = |T(D)|$ となり、 D は T によって情報の損失を伴うことなく $T(D)$ に圧縮されたことになる。

フィルターとしての述語 P は正確に $D = L(G(D))$ となるように $D(G)$ を推論すれば不要であるが、 D は有限であるから、 $|G(D)|$ の値はあまり小さくならないことが多い。(上の D を有限オートマトンで表現すると状態は最少の場合でも 103 個必要であり、それを 2 進数でコード化すると約 600 のスペースが必要であり、推移表の面積は $600 \times 2 = 1200$ となる)。このことと、 D がある母集団からのサンプルである場合には、推論問題は付帯条件を除いた (2.4) だけで扱う方が自然であるので、先ず (2.4) を考え、更に述語 P を考えるのである。

3. 圧縮の可能性

圧縮の可能性、度合の問題はデータの複雑性の問題として

扱うことができる。データ D の複雑さ, すなわち, 可能な圧縮の割合は

$$H(D) = \min\{|Q| + |P|; L(Q, P) = D\} \quad (3.1)$$

で定義できる。定義から直ちに

$$H(D) \leq |D| \quad (3.2)$$

であることが分る。

多くの D に対して $H(D)$ の値が十分小さくなることを期待したいのであるが, 悲観的材料として Kolmogorov の複雑さに関する結果がある。即ち, $x \in \Sigma^*$ の複雑さを $K(x)$ とすると,

$$K(x) \leq |x| \quad (3.3)$$

$$\#\{x; K(x) < l - m\} / \#\{x; |x| = l\} \leq 2^{-m+1} \quad (3.4)$$

であることが知られている。つまり, 大多数の x に対して, 情報の損失の多い大幅な圧縮は期待できないことに分る。しかも, D と 適当に順序付けられた系列であると看之れば, 一般に

$$K(D) \leq H(D) \quad (3.5)$$

であるから, 任意のデータを文法推論による圧縮法により, 大幅に圧縮することは不可能である。

しかし, 集合 D の要素の順序付けの仕方によって $H(D)$ の値が小さくなる可能性はある。もし, D に明確な規則性がある場合, あまは D についての付加的情報が与えられている場合には大幅な圧縮も可能である。一方, 再び Kolmogorov の理論から, $K(x)$ の値が大

至ければ大至り程又は random に近いことが知られているので、 $H(D)$ の値を下げることが不可能である場合には、 D は random と考えて差しつかえなく、その分布も確定するので、 D について考えられる統計的問題はすべて解決されることになる。この場合は、その分布を記録するだけで十分であるので、再び大幅な圧縮への道が開ける。勿論この意味での情報の損失は伴うことになるが、 D に関しては、 D そのもの以外に、いかなる決定論的で有効な特徴抽出の手法もないのである。

4. 右延長文法とその推論

文法推論において、データの特徴を記述するシステムは、従来のオートマトン、文法に限定されたものでは無い。ここでは、推論、特徴抽出の最も素朴な形式の定式化とみなされる文法種を導入し、1つの推論を考へることになる。

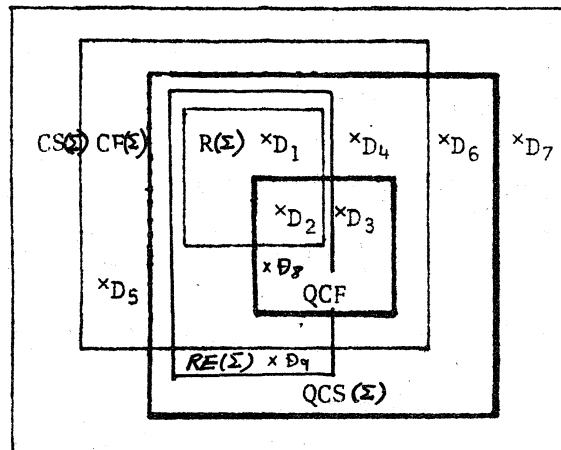
4.1. 擬文法とその性質

文脈依存文法 $G = (V, V_T, P, S)$ による *sentential forms* の全体のなる集合 $S(G) = \{w; S \xrightarrow{+} w\}$ を G において生成される擬文脈依存言語 (qcsL) といふ。 $\xrightarrow{+}$ は1回以上のルール適用による導出を示す。qcsL について考へる際には、 V_T を明示する必要はないから、 $G = (V, P, S)$ といふ、これを qcsG と呼ぶ。 $G = (V, V_T, P, S)$ が文脈自由であるといは、qcfG といふ。

gcsq $G=(V, P, S)$ において, P の要素が $S \rightarrow \alpha, \alpha \rightarrow \alpha\beta$ の型であるとき, G を右延長文法といい, 以後 $S := \{\alpha; S \rightarrow \alpha \in P\}$, $P := \{(\alpha, \beta); \alpha \rightarrow \alpha\beta \in P\}$ とする.

記号の集合 Σ 上の csl の全体を $CS(\Sigma)$, cfl の全体を $CF(\Sigma)$, 正則集合の全体を $R(\Sigma)$, gcsl の全体を $QCS(\Sigma)$, gcf の全体を $QCF(\Sigma)$, 右延長言語の全体を $RE(\Sigma)$ とする.

$\Sigma \geq 2$ ならば, この言語族には下図のような関係がある:



$$\begin{aligned}
 D_1 &= \{A\} \cup \{A^{2n+2}; n \geq 0\}, & D_2 &= \{A\}^* \\
 D_3 &= \{A^n B A^n; n \geq 0\}, & D_4 &= \{A^n B^n; n \geq 0\} \\
 D_5 &= \{A B^n A B^{n+1}; n \geq 1\}^*, & D_6 &= C(S(G)); L(G) = \{a^{2^n}; n \geq 0\}, \\
 D_7 &= \{A^{2^n}; n \geq 0\}, & D_8 &: \{A, B\} \text{ 上の Dyck 言語,}
 \end{aligned}$$

ここに C は G の アルファベットの $\{a, S, A_1, \dots, A_p\}$ とおくと,

$$\begin{aligned}
 C(a) &= A^{p+1} B, & C(A_i) &= A^{p-i+1} B^{i+1}, & C(S) &= S \\
 & & & (1 \leq i \leq p)
 \end{aligned}$$

で定義される homomorphism である。 D_9 については [14] を参照してほしい。

なお、 $\#\Sigma = 1$ の場合には、この言語族間には、

$$QCF(\Sigma) \subseteq R(\Sigma) = RE(\Sigma) = QCS(\Sigma) \subseteq CS(\Sigma) \quad (4.1)$$

の関係がある。QCF(Σ)等の他の性質について少し述べる。

定理 1. $\#\Sigma \geq 2$ とする。QCS(Σ), QCF(Σ), RE(Σ) は "どの" の演算の下でも閉じている: (1) union, (2) product, (3) *, (4) intersection, (5) Σ^* に関する complementation, (6) homomorphism.

また、QCS(Σ), QCF(Σ) は 正則集合 との intersection の下でも閉じている。

定理 2. $\Sigma = \{A\}$ とする。

(1) $D \subseteq \Sigma^*$ が bcf l $\Leftrightarrow D = \delta(D) \cup X \cdot \{A^r\}^*$ とする有限集合 $X \subseteq \Sigma^*$ と整数 r が存在する。

(2) $D \subseteq \Sigma^*$ が gcs l $\Leftrightarrow D = F \cup X \cdot \{A^r\}^*$ とする有限集合 $X, F \subseteq \Sigma^*$ と整数 r が存在する。

(3) QCF(Σ) は union の下で "必ず" 閉じているが, $D_1, D_2 \in QCF(\Sigma)$ が "無限" 個あるならば, $D_1 \cup D_2 \in QCF(\Sigma)$ である。

定理 3. $QCS(\Sigma, k) = \{D \subseteq \Sigma^*; (\exists \text{ gcs } G)[S(G) = D \ \& \ o(G) = k]\}$,

$$o(G) = \max\{|x|; (\exists p)[\alpha \rightarrow \beta \in p]\}$$

とすると, (1) $QCF(\Sigma) = QCS(\Sigma, 1)$, (2) $QCS(\Sigma, k) \subseteq QCS(\Sigma, k+1)$,

(3) $QCS(\Sigma) = \bigcup_{k=1}^{\infty} QCS(\Sigma, k)$ であり, Σ は

(4) $D \in QCS(\{A\}, k) \Leftrightarrow D = F \cup X \cdot \{A^r\}^*$ とする有限集合 $X, F \subseteq \Sigma^*$ と整数 r が存在し, 各 $z \in F$ に対して $|z| < k$ とする。

4.2. 右延長文法の推論

前節でおなじみのように擬文法による言語族は通常の言語族とは極めて異なった性質をもつが、文法は文の生成のために何も補助的な記号、状態等が必要としないので推論には都合がよい。以下に右延長文法の推論について簡単にのべる。

長さの短かい順に並んだ n 個のサンプル系列の集合

$$D = \{x_1, x_2, \dots, x_n\}; |x_i| \leq |x_{i+1}| \quad (1 \leq i \leq n) \quad (4.2)$$

が与えられたとき条件

$$D \subseteq S(G) \quad \& \quad S(G)|_D = D \quad (4.3)$$

を満たすような右延長文法 G を推論する方法を考へる。ここに

$$M = \max \{i; i \geq 1 \& |x_i| < |x_n|\} \quad (4.4)$$

とすると

$$S(G)|_D = \{x \in S(G); |x| \leq |x_M|\} \cup \{x_i; M < i \leq n\} \cap S(G) \quad (4.5)$$

である。推論の手順の詳細は省略するが、本質的な部分は、各 x_i に対して x_i を導出するような x_j ($j < i$) を見つけることにある。 x_j が見つければ、 $x_j \Rightarrow x_i$ を導く規則を検討して、 D の長さまでの語で D に属する語を導出しなければ、その規則を P の元として採用する。 $x_j \Rightarrow x_i$ を導くすべての規則が D 以外の $|x_n|$ 以下の語を導出するが、 x_j が見つかるときは、 x_i を S の元として採用する。このようにして、右延長文法 $G = (\Sigma, P, S)$ を決定する。最後にこの方法による推論過程の状況を示しておこう。

i	x_i	C	S	P
1	011		011	
2	00111	(0,01)	00111	
3	01101	(011,01), (01,10)	01101	
4	0010111	(<u>00</u> ,10), (0,0101), (0,01), (001,01)		(00,10)
5	0011101	(<u>0011</u> ,01), (0,01)		(111,01)
6	0110101	(01101,01), (<u>0110</u> ,10), (011,0101), (01,1010), (01,10), (011,01)		(10,10)
7	001010111			
8	001011101			
9	001110101			
10	011010101			
11	00101010111			

5. おわりに

文法推論による情報圧縮について基礎的考察を行い、その可能性について論じた。また、右延長文法とその周辺の文法を考へ、推論についての簡単考察を試みた。これらの中には筆者のこの方面での研究の方針を述べただけのところも少くない。したがって、解決可能な多くの問題が残されている。データの是認、目標と呼んで来たものの是式化、変換の是式化、付加の情報の扱い、推論手順の比較、改良などの問題がある。これらについては別の機会に論じるつもりである。

参考文献

1. Moore, E.F. Gedankenexperiments on sequential machines, in Automata Studies, Princeton Univ. Press (1956).
2. Solomonoff, R.J. A formal theory of inductive inference, Inform. Cont. 7 (1964), 1-22.
3. Gold, E. M. Limiting recursion, J. Symbolic logic 30 (1965), 28-48.
4. 榎本肇: カンフォルム記号列によるオートマトンの記述と構成, 日科技連 数学計画シンポジウム「オートマトン」(1971) 141-163
5. Feldman, J. Some decidability results on grammatical inference and complexity, Inform. Cont. 20 (1972), 244-262.
6. Biermann, A. W. An interactive finite-state language learner, 1st US-Japan Comp. Conf. Proc. (1972), 13-20.
7. Banerji, R. B. Theory of problem solving, Amer. Elsevier (1969).
8. Smullyan, R. Theory of formal systems, Princeton Univ. Press (1961).
9. Tanatsugu, K. On inference of a concise description of finite and infinite concepts, RIFIS-RR 34 (1973), 1-13.
10. Zvonkin, A. K. & Levin, L. A. The complexity of finite objects and the development of the concepts of information and randomness by means of the theory of algorithms, Rus. Math. Surv. 25 (1970), 83-124.
11. Tanatsugu, K. & Arikawa, S. On characteristic sets and degrees of finite automata, RIFIS-RR 45 (1974), 1-14.
12. Arikawa, S. On the languages defined by sentential forms of context-free grammars, RIFIS-RR 17 (1970), 1-12.
13. Arikawa, S. Closure and nonclosure properties of quasi context-sensitive languages, RIFIS-RR 37 (1973), 1-25.
14. 大芝猛, 有川: 右辺文法による言語の族の大きさについて 京大数研講義録 213 (1974) 104-122.