

構文情報を用いた一つの音声認識法について
——未知語の取り扱い——

九州大学理学部	武谷 峻一
九州大学工学部	河口 英二
九州大学工学部	田町 常夫

1. まえがき

連続音声の認識における言語的な情報の重要性は古くから指摘されており、近年いろいろなレベルの言語情報を利用する認識システムが研究されている^{(1),(2),(3)}。これらの多くは、入力文に対して文法的な制限を設け(対象を特定の文法構造をもつ言語分野に限定するなど)、その文法規則を利用することにより認識の質を高め、処理効率の向上を図っている。こうしたシステムでは、処理にあたっていざれも“単語”という概念を媒介としなければならず、単語辞書が不可欠である。

辞書を用いる場合、認識の対象を拡大しようとするときそれに伴って内蔵すべき語彙の増大は避けられない。しかし、実

際に入力に現われる単語の頻度にはいろいろな偏りが見られる⁽⁴⁾。天気予報文448文章を対象としてそれに用いられている単語151語の使用頻度を調べた結果、最も頻度の高い13語で全体の50%、150語では97%以上が尽されている。そして残りの2~3%のために頻度1または2の単語約100語が必要となっている。こうした傾向は自然語の一般的法則としてよく知られている。したがって、辞書の利用効率から見れば、すべての単語を辞書に登録しておくよりも、頻度の低い単語(たとえば固有名詞など)を“未知語”として検出できる方式であることが望ましい。

こうした観点から、本稿では文の構文的な構造を利用して未知語を検出する一つの方法を提案する。

2. 未知語取り扱い上の問題点

未知語を扱う場合、入力系列を単語に区切ってゆく操作が必要である。単語への分割操作ということは通常の言語処理ではあまり問題にならない。あるいは、入力があいまいさのない文字列であれば、見出し語と完全に一致する部分を単語として区切れればよい。この場合、未知語の検出は比較的容易である。とくに英語では空白によって単語が区切られているのでほとんど問題がない。しかし音声認識では、一般に入力系列は物理的な特徴によって識別されたおとの不完全な(あ

いまのたの記号列であり、見出し語との照合において完全一致となるものを探し出せることは期待できない。そこで、相対的な一致の尺度を設定して判定せざるを得ないが、そのために分割の多様性が生じることになる。したがって、未知語を取り扱うには、何らかの形で構文的な情報と利用する分割法が必要となる。その際、一般によく行われる左から右への分割法は不適當であろう。

3. 未知語を考慮した単語分割

2.で述べた点を考慮して、構文を利用し未知語の検出およびその構文上の機能の推定が可能な分割法について述べる。

3.1 入力に対する仮定

まず、入力に対してつぎのような仮定を設ける。すなわち、対象とする音声は以下に定める生成文法 G により生成された文を区切り記号の部分で区切りながら連続音声として発話したものであると仮定する。

文法 G は

H : 開始記号,

$S = \{a_1, a_2, \dots, / \}$: 音韻記号と区切り記号の有限集合,

$W = \{w_1, w_2, \dots, w_N, w_{N+1}, \dots\}$: 単語の集合, N は辞書に登録されている単語数,

$C = \{c_1, c_2, \dots\}$: 構文的機能(品詞, 区切りなど)を表わ

す記号の有限集合,

$V = \{\alpha, \beta, \dots\}$: 書き換えのための変数 の有限集合

とする. S は終端記号の集合, $W \cup C \cup V$ は非終端記号の集合であり, 生成規則 P はつぎの4つの形式からなる.

$$(i) \quad H \rightarrow / \alpha /$$

$$(ii) \quad C_i \beta C_j \rightarrow C_i \gamma C_k \delta C_j$$

$$(iii) \quad C_i \rightarrow w_n$$

$$(iv) \quad w_n \rightarrow A_1 A_2 \dots A_k$$

ただし, 区切り記号は単語でもあり構文的機能を表わす記号でもある ($/ \in W, / \in C$). また, $\epsilon \in V$ ($|\epsilon| = 0$) とする. G の文は一般に

$$Q = / A_1 A_2 \dots A_k / A_{k+1} \dots / \dots A_m /$$

となる.

3.2 単語分割のアルゴリズム

Q を連続発声した音声は, その物理的な特徴によって識別され, 一たん音韻記号列 Q' ($\in S^*$) に変換される. この記号列は一般に誤って識別された音韻記号 w_i を含む不完全な記号列である. 本稿ではこの Q' を入力系列として取り扱う. Q' 中の音韻記号列と辞書 $D = \{w_1, w_2, \dots, w_n\}$ ($\subset W$) 中の単語との相対的な一致の尺度はつぎのように定める. いま, $S_0 = S - \{/\}$, $x \in S_0^*$ とするとき, x と $w \in D$ との一致の度を

$\omega(w, x)$ とする. このとき, w に固有のしきい値 $\theta(w)$ を考え, $\omega(w, x) \geq \theta(w)$ であれば w と x が一致したとする.

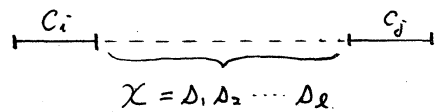
G の文 Q を認識する立場から見ると, 未知語 $(w_{N+1}, w_{N+2}, \dots)$ に関しては G の生成規則 (iii), (iv) が欠落してゐたと見做すことができる. そこで単語分割にあたっては, 生成規則 (i), (ii) の形式を構文的な情報として利用する. すなわち, 関数 $\pi: C \times C \rightarrow 2^C$ を考え, 左および右側の構文的機能がそれぞれ C_L, C_R であるとき, $C_L - C_R$ 間の構文的機能を $\pi(C_L, C_R)$ で規定する. この π と辞書 D を用いる分割のアルゴリズムをつぎのように定める (例参照).

分割は段階的に行なうこととし,

いま注目している音韻記号

列 $x = a_1 a_2 \dots a_l$ ($x \in S_0^*$, Q の

部分記号列) に関して, $C_L = C_i, C_R = C_j$ とするとき,



(I) $\pi(C_i, C_j) = \emptyset$ の場合.

$C_i - C_j$ 間に許される構文的機能がないうちの場合で, 前段での C_i あるいは C_j の決定が誤りである. 前段へ戻り新たに別の C_i あるいは C_j を選ぶ.

(II) $\pi(C_i, C_j) \neq \emptyset$ で, y が x の部分記号列 y ($y \prec x$ と表わす) に対して

$$T_y^x = \{w \in D \mid \omega(w, y) \geq \theta(w), w \in \psi(c), c \in \pi(C_i, C_j)\}$$

(ただし, ϕ は $C \rightarrow 2^D$ であり $\phi(C)$ で構文的機能 C をもつ単語の集合を表わす) としこれとき, すべての $y \prec x$ に対して $T_y^x = \phi$ の場合.

すなわち, $C_i - C_j$ 間に許される構文的機能はあるが, x 中に $C \in \pi(C_i, C_j)$ をもつ既知語 $w \in D$ と一致する部分がない場合は, x を未知語と決定する.

(III) $\pi(C_i, C_j) \neq \phi$ で, ϕ が $T_y^x \neq \phi$ となる $y \prec x$ が存在する場合.

すなわち, $C_i - C_j$ 間にある $C \in \pi(C_i, C_j)$ をもつ既知語 $w \in D$ と一致する部分がある場合は,

$$Y_x = \{y \prec x \mid T_y^x \neq \phi\},$$

$$\omega(w_0, y_0) = \max_{y \in Y_x} \max_{w \in T_y^x} \omega(w, y)$$

とおき, π の y_0 の部分と w_0 とする. そして Q' の y_0 に相当する部分と w_0 の構文的機能が置き換えて次段へ進む. ところが, 次段が (I) の場合で逆戻りしたときは, 順次 ω , ω の番目に大きい $\omega(w, y)$ を選ぶことになる.

以上の分割法により, 未知語部分は最後に音韻記号列のまゝ残る. いま, π の両隣の構文的機能が $C_L = C_i, C_R = C_j$ であるとすれば, その未知語の構文的機能は $\pi(C_i, C_j)$ のいずれかであると推定できる.

実際の命割操作では, $c \in \mathfrak{K}(C_L, C_R)$ に対して, C_L に隣接, C_R に隣接, 両方に隣接, 中間位置などの情報が利用できる. しるが, Z , $W(W, X)$ の計算が単語・語尾からもできる形式であれば, これらの情報により照合の回数を減らすことが可能である.

3.3 \mathfrak{K} の構成

上に述べたアルゴリズムを用いて Q' を単語に命割してゆくとき, あらかじめ \mathfrak{K} をどのように定めておくかが問題となる. \mathfrak{K} としては, まず

$$\mathfrak{K}_1(C_i, C_j) = \{c \in C \mid c_i \alpha c_j \rightarrow c_i \beta c \gamma c_j \in P\}$$

が考えられる. これは単語命割を G の生成過程と全く同じに追わせることを表わしている. この場合, たとえば Q の両端の区切り記号 $/$ に関して P に $/ \alpha / \rightarrow / \beta c \gamma /$ がただ1個しかたゞときには, $\mathfrak{K}(/, /) = \{c\}$ となる. この c に相当する部命が未知語であれば, 命割は才1段目で未知語を検出して止まり, Q' 全体が未知語と決定される.

つぎに, \mathfrak{K} として $L(G)$ のすべての文に対して任意の $C_L - C_R$ 間の構文的機能調べ, 可能な構文的機能すべて $\mathfrak{K}_2(C_L, C_R)$ とする方法がある. 一般に, 任意の C_L, C_R に対して $\mathfrak{K}_1(C_L, C_R) \subseteq \mathfrak{K}_2(C_L, C_R)$ となる (例参照). この \mathfrak{K}_2 では, Q の両端の $/$ に対して $\mathfrak{K}(/, /) = C$ (すべての構文的機能が

含まれる)となり、分割の才1段目でD内のすべての単語とQのすべての部分系列との照合が必要となる。

以上の点を考えれば、至₁は未知語が存在するQ'に対しては不適當であり、至₂は単語との照合回数が多く処理効率が悪く、したがって、現実の至として至₂を採用し、至(C₁, C₂)内のC'に対して優先度を設けて処理を行おうが、あるいは未知語を許す構文的機能を限定しておく(たとえば、助詞、助動詞などは未知語を許さず、固有名詞は許すなど)などの制限を設ける必要があるだろう。至₂、いずれの方法を用いても、3.2のアルゴリズムではQ₁とL(G)を受け入れる可能性があり、単語分割のあとに構文的なチェック(未知語に対しては推定された構文的機能をあてはめてチェックを行おう)が必要である。

4. おまじ

本稿で示した分割法の特徴は、構文的な情報を形式的に至₂で取り扱い、上から下の分割してゆくことにより、未知語の検出およびその構文上の機能の推定が可能となる点である。しかし、本方式で未知語として処理される音韻記号列には、(1) 真の未知語である場合と、(2) 既知語でありながら音声の物理的な特徴による識別段階で誤りが大きくなり、分割を誤った場合とがある。この両者を区別する情報はなく、いずれも

同一の処理を行わねばならぬ。また、連続して未知語は取り扱えないなどの問題点が残されている。

現在、本方式による認識実験を行っており、近々その成果を報告する予定である。

参考文献

- (1) A. Newell et al. : "Speech Understanding Systems, Final Report of a Study Group", Carnegie-Mellon Univ., Pittsburgh (1971-05)
- (2) 斎藤他 : "音声認識における言語情報の利用", 昭和49年電気四学会連合大会論文集, No. 193 ~ 199 (昭49-10)
- (3) 武谷, 河口 : "言語構造情報と利用した連続音声認識システムのシミュレーション", 電子通信学会論文誌(A), Vol. 56-A, No. 9 (昭48-09)
- (4) 河口, 武谷, 田町 : "音声認識における文法規則の利用", 昭和49年電気四学会連合大会論文集, No. 195 (昭49-10)

構文的機能と単語

T: 明日は, 明後日は, ... C: 曇り, 晴れ, 雨, ...

F: のち, 一時, 時々, ... X: です, でしょう, ...

生成規則 (i), (ii) の形式のみ)

$$H \rightarrow / \alpha /$$

$$/ \beta C \rightarrow / T \delta C$$

$$/ \alpha / \rightarrow / \beta C /$$

$$T \delta C \rightarrow T C \epsilon C$$

$$/ \alpha / \rightarrow / \beta C \gamma /$$

$$C \epsilon C \rightarrow C F C$$

$$/ \beta C \rightarrow / T C$$

$$C \gamma / \rightarrow C X /$$

文

$$/ T [C F] C [X] / \quad [] \text{は省略可}$$

{ 明日は晴れ。 明後日は雨でしょう。
明日は曇りのち晴れです。 ... }

Φ_1

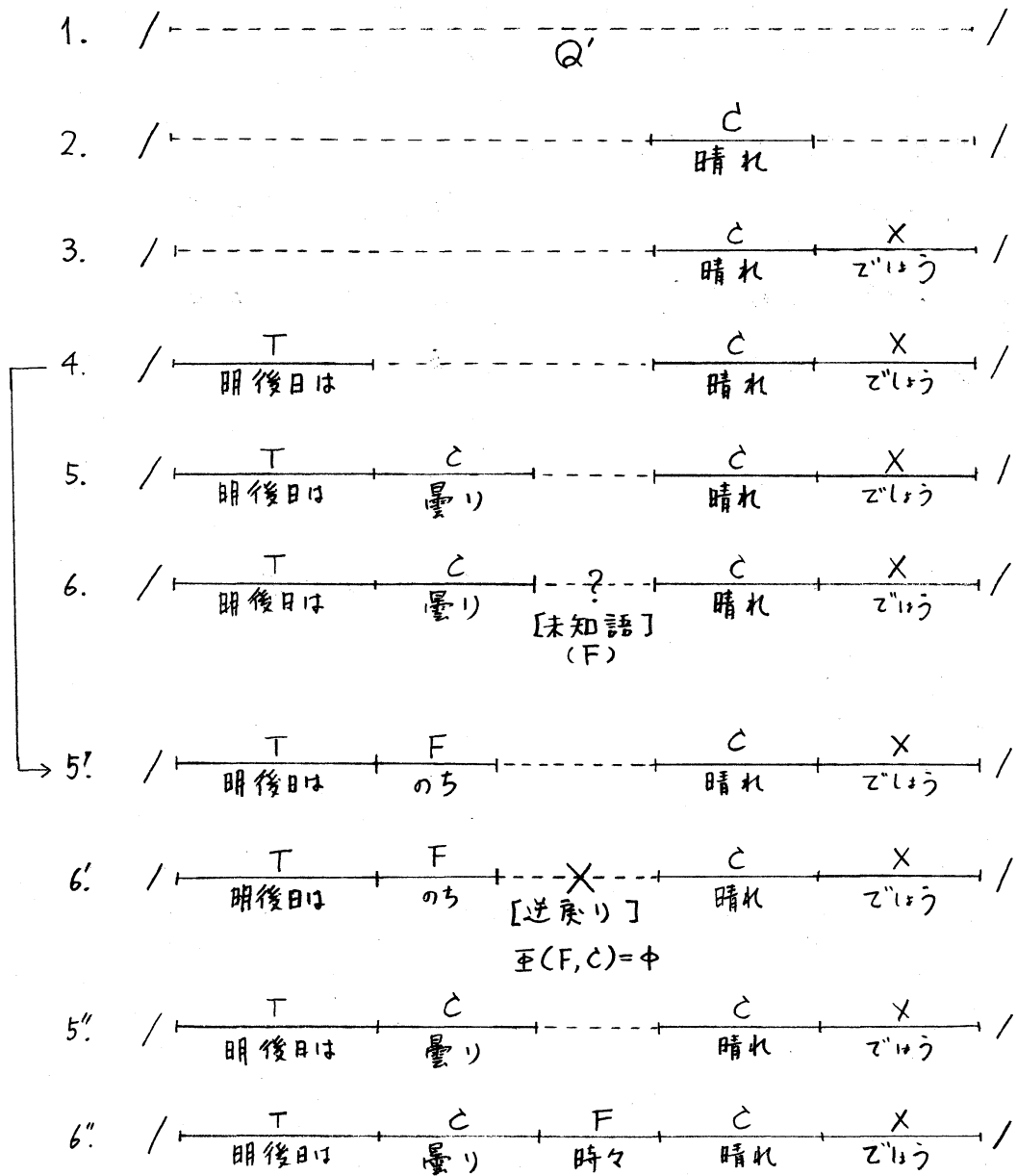
Φ_2

$\begin{matrix} C_R \\ C_L \end{matrix}$	/	T	C	F	X
/	{C}	φ	{T}	φ	φ
T	φ	φ	{C}	φ	φ
C	{X}	φ	{F}	φ	φ
F	φ	φ	φ	φ	φ
X	φ	φ	φ	φ	φ

$\begin{matrix} C_R \\ C_L \end{matrix}$	/	T	C	F	X
/	{T, C, F, X}	φ	{T, C, F}	{T, C}	{T, C, F}
T	{C, F, X}	φ	{C, F}	{C}	{C, F}
C	{C, F, X}	φ	{F}	φ	{C, F}
F	{C, X}	φ	φ	φ	{C}
X	φ	φ	φ	φ	φ

例

単語分割の進み方



例(続き)