

# AP電経済ニュースからの定型パターンの抽出

## 6E-4

浦谷 則好      加藤 直人      相沢 輝昭

NHK放送技術研究所

### 1 はじめに

これまでに報告しているように放送用の英日機械翻訳システムの研究を進めている。<sup>1)</sup>開発中のシステムは1989年8月以降衛星第1放送において日本語テロップの作成に試用されている。翻訳の精度の向上のために、辞書・文法の改良のほか、定型的なパターンを直接日本語に変換する方式との結合(ハイブリッド方式)も検討している。このために、AP電の英語ニュースから定型パターンの抽出を行った。

### 2 ハイブリッド方式

英語ニュース文には、例えば

In Kuala Lumpur, Malaysian tin closed at 28.43 dollars per kilo, down 3 cents.

などの、定型的なパターンがよく出現する。このような文は解析-変換-生成という手順を踏まなくても、「クアラルンプールでマレーシアのすずは1キロ当り、#Aドルで引けた。#Bセントの#Cである。」という、日本語パターンを用意しておけば、#Aに28.43、#Bに3、#Cに「下げ」を(upなら「上げ」)代入すれば、高品質の日本語訳文を高速に得ることができる。もちろん、Kuala Lumpur, kilo, Malaysian tinも変数として考えることができる。AP電のデータをもとに、このような定型パターンの利用の検討を進めている。上例のように文全体がパターンとなり得ることもあるが、一般的にはパターンは文の一部として現われると考えられる。我々は、解析-変換-生成を基本にし、文中のパターンは一括して捉え、解析のあいまいさを減らし、高速化と訳文品質の向上を図っていく予定である。これを、ハイブリッド方式と呼ぶことにする。この考え方の一部はすでに、固有名詞や数量表現の処理に利用している。<sup>2)</sup>

### 3 AP電経済記事に出現する定型パターンの抽出

AP電1年分('89.4~'90.3:ただし、データの欠落等のため実質約300日分)の経済ニュースからパターンを抽出するために、まず単語の連続

出現頻度を調べた。長いパターンを抽出するにはできる限り長い単語列の頻度を取った方が良いが、記憶容量と計算時間を考慮して、10語連続を調べたことにした。テキストから1語づつずらして全ての10語連続の頻度を求めた。こうしておけば、9語連続以下の頻度情報も求めることができるし、これくらいの長さがあれば、11語以上の頻度も後述する方法で推定することができる。ただし、文末を越えての語の連続は認めないこととし、コンマや括弧等も1語と見なすことにした。さらに、重複記事と判別できるものは統計対象から除外した。結果として得られた頻度表を表1に示す。一番出現している"futures trading on ..."は612回(全体(145万)の0.04%)で、1日に平均2回程度出現していることになる。

表1. 10語連続の出現頻度(経済記事) (上位10位)

futures trading on the New York Coffee, Sugar and trading on the New York Coffee, Sugar and Cocoa	612
on the New York Coffee, Sugar and Cocoa Exchange	611
a bid price of Y dollars a troy ounce,	608
(dollars per metric ton) cash Y-Y (Y-Y	349
at a bid price of Y dollars a troy ounce	345
in London at a bid price of Y dollars a	298
London at a bid price of Y dollars a troy	264
, with Y up, Y down and Y unchanged	264
and dealer buying-dealer selling U.S. dollars at Y-Y	258
Hong Kong	247

[Y: 数字列]

#### 3.1 10語連続頻度表からの長短頻度表の作成

表1からも分かるように、10語連続は実はもっと長いパターンの部分である可能性がある。(上位3位は12語以上のパターンの一部と考えられる。)また、10語ではパターンをなさないが、9語以下でパターンとなっているものも考えられる。パターン表現としてはできるだけ長いものを抽出したい。また、曜日表現など変数として扱うべきものを個別に扱っては、パターンが増えるだけで有効なパターンが現われてこない。そこで、同一と見なした方が良いと思われる単語群を、曜日を全てWednesday, 月は通常表現はAug., フルスベルのものはNovember, 記号的表現はAprにするなどの同一化を行うことにした。短いパターンの抽出のために、10語連続頻度表から後ろの1語を順次除いてゆき、9語・8語・7語・6語の頻度表

を作成した。また、長いパターンを求めるために、11語、12語、・・・の頻度表を以下のように作成した。

#### < n + 1 語頻度表の作成 >

n語のある単語列 $\alpha$  (例えばabcdefghij)の後n-1語と別の単語列 $\beta$  (bcdefghijk)の前n-1語が一致したときは、一致した部分(bcdefghij)を含むn+1語の単語列 $\gamma$  (abcdefghijk)を生成する。この単語列の頻度は $\alpha$ 、 $\beta$ の内小さい方の値とする。ただし、 $\alpha$ 、 $\beta$ の頻度の差が大きいとき(片側が2倍以上)にはn+1語単語列を生成しない。

#### 3. 2 定型パターンの自動抽出

あるn語の単語列の頻度がmであれば、前後どちらかの語を除いたn-1語の単語列の頻度は必ずm以上となるから、単純に頻度が大きいからといってそれをパターンと考えるわけにはいかない。大量のデータを扱うときには、真のパターンとなり得るものを何らかの基準をもとに自動的に抽出したい。これには、仕事量の基準<sup>3)</sup>を採用することにした。すなわち、単語列 $\alpha$ の総仕事量を

$$W(\alpha) = |\alpha| \times n(\alpha)$$

(ここで $|\alpha|$ は単語列の長さ(語数)

$n(\alpha)$ は単語列の出現頻度)

で定義する。(ただし、文献3と異なって $|\alpha|$ を用いて、仕事量の差でなく、仕事量そのものを用いている。)そして、単語列 $\alpha$ があるとき、この先頭の単語を除いた単語列 $\beta$ とし、最後尾の単語を除いた単語列 $\gamma$ としたとき、

$$W(\alpha) > W(\beta) \quad \text{かつ} \quad W(\alpha) > W(\gamma)$$

のとき、 $\alpha$ を有効なパターンとして認定する。

ただし、単語列 $\alpha$ を部分列として含む1だけ長い単語列のうち最大頻度を持つものを単語列 $\delta$ としたとき、 $\alpha$ の頻度を $n(\alpha) - n(\delta)$ で更新し、これが2以下になるものは除外する。つまり、 $\alpha$ は $\delta$ で表現されないものの代表と考えるのである。

3. 1で得られたデータをもとに、上記の基準で有効パターンの抽出を44語連続まで行った。11語連続以上で頻度200(0.014%)以上のものを表2に示す。10語連続以下の上位3位までを表3に示す。

#### 4 おわりに

表2、表3を見ると表1ではノイズ中に埋もれていたパターンの候補が顕現していることが分かる。しかし、単語の同一化処理の不足のために、本来はもっと長いと推測されるパターンの一部が現われていると思われるものもある。ところで、n語の連続頻度表からn-1語が一致するものを取ってきて、異なっている

表2. 11語以上の出現パターン(経済記事) (頻度200以上)

12語	-----	
	futures trading on the New York Coffee, Sugar and Cocoa Exchange	608
	to close at Y U.S. dollars, compared to Wednesday's Y	240
13語	-----	
	's foreign currency prices (dealer buying-dealer selling) quoted in U.S. dollar	211
	the U.S. dollar closed at Y yen on the Tokyo Foreign Exchange Market	208
	In Kuala Lumpur, Malaysian tin closed at Y dollars per kilo.	208

表3. 10語以下の出現パターン(経済記事) (上位3位)

10語	-----	
	the New York Coffee, Sugar and Cocoa Exchange Wednesday (dollars per metric ton) cash Y-Y (Y-Y	609
	, to close at Y U.S. dollars, compared to	345
		244
9語	-----	
	European Currency Unit, a basket of European currencies	205
	Futures trading on the Chicago Board of Trade Wednesday of Y dollars a troy ounce, up from	193
		184
8語	-----	
	--Y West German marks, up from Y	358
	Y futures on the N.Y. Cotton Exchange Wednesday	205
	Y dollars a troy ounce, up from	185
7語	-----	
	--Y French francs, up from Y	381
	--Y Swiss francs, up from Y	376
	--Y Dutch guilders, up from Y	373

語を調べれば、同一化すべき単語(変数)の候補も得られると考えられる。この結果は別の機会に報告したいと考えている。

ここでは、経済記事についての結果を述べたが、一般記事についても10語連続頻度を取ってみた。結果は一番頻度の高いものでも116回(全体(1600万)の0.0007%)であった。頻度50以上のものは8個で、一般記事にはそんなに長いパターンは現われないことが判明した。

今後は、得られたパターンの候補を精査し、各パターン処理方法を決めていく予定である。必要とあれば、節・句・品詞などの解析単位の見直しも行う予定である。さらに、スポーツ記事や一般記事(の短いパターン)についても、解析を進めていく予定である。

#### 【参考文献】

- 1) 浦谷ほか「英語ニュースの機械翻訳」情報処理研究会NL78-18(1990)
- 2) 加藤ほか「英日機械翻訳における固有名詞処理」情報処理学会第40回全国大会2F-2(1990)
- 3) 北ほか「テキスト・データベースからの慣用表現の自動抽出」ATR Technical Report TR-1-0027