

# 動的計画法による漢字仮名混り文の単位切りと仮名ふり

藤崎哲也助 (日本アイ・ビー・エム(株) 東京サイエンティフィック・センター)

## 1. はじめに

日本語の言語処理が大学や会社の研究室の範囲を出て、様々な分野で実用的にその技術が用いられるようになってきたが、そこで依然として重要な基本的問題として日本文の入力の問題と漢字仮名変換の問題がある。

この漢字仮名変換とは通常の日本語の表記法に従いバタ書<sup>1</sup>または漢字仮名混り文を入力として、わかち書き<sup>2</sup>または仮名文を出力するものであり、バタ書<sup>1</sup>を適切な単位でわかち書き<sup>2</sup>単位切り<sup>3</sup>の問題と、わかち書き<sup>2</sup>の漢字の読みを文脈より決定する仮名ふりの問題の2つの問題を含んでいる。

この漢字仮名変換の重要性については特に言及するまでもないが、文書よりの索引(KWIC/KWOC)の自動作成、文献検索のための自動キーワード切り出し<sup>4</sup>などを実現するために不可欠の技術であり、また文書の点字化、将来の文書の自動朗読<sup>5</sup>などにも結びつく技術である。

日本アイ・ビー・エム(株) 東京サイエンティフィック・センターでも、従来より将来のオフィス・オートメーション実現のための中核として日本語文書処理システム「ことばま」の研究開発を行って来ているが、<sup>[1,2,3,4]</sup> ところで、既存の漢字<sup>1</sup>入力<sup>2</sup>または文書<sup>3</sup>に対して「ことばま」の索引自動作成、キーワード文献検索などの文書処理サービスを利用可能とするために漢字仮名変換の研究を行って来ている。

過去においても、この漢字仮名変換には、いくつかの試みが報告されている。<sup>[5,6,7]</sup> また、その目的により、特に単位切りにおいては、文節の単位

の単位切りを目指すもの、自立語/付属語の単位切りを目指すもの、さらに自立語が複合語である時に語基の単位切りまで目指すものなどの差異があり、それらを一律に比較することはできない。

文節の単位までの単位切りにおいては、字種の変り目(ひらがなからひらがな以外への変り目)を文節の始まりとするという単純な規則が有効であることがよく知られており、この単純な規則<sup>8</sup>だけにより84%の分割精度を得たとの報告が行われている。<sup>[8]</sup> 単位切りの単位が短くるとさらに難しくなる。漢語が多く2文字単位で構成されることや、複合語分解の分布を利用し規則と若干の辞書用漢字リストで分割を試みる例もあるが、<sup>[9]</sup> 最終的には語基の辞書を整備し、分割の際に教員<sup>9</sup>の辞書を参照しなければならない現状である。

一先、仮名ふりにおいても、単位切りと同様に、両隣の文字の字種が漢字以外なら訓読み、その他の場合は音読みとするというような単純な音訓規則を利用すると、ある程度の漢字は読みが一通りしかないのである。<sup>[7]</sup> 85%程度の精度を得ることはできる。<sup>[7]</sup>

このように、文節切り、仮名ふり共に簡単な規則で85%程度の精度は容易に達成できるのだが、単位切りの単位をより細くする、精度を高め、などのためにこれらの基本規則の例外を辞書もしくは例外コードとして蓄積することになる。自然言語は奥の深いものである。このようなアプローチは最終的に避けられないと考えられるが、基本的規則の与える精度が低ければ低い程、辞書などの形で蓄積しなければならぬ

い情報の量が増大となり、かつ分野への依存度が高まるのは避けられない。

筆者は、従来より文字認識・音声認識などのパターン認識の分野で用いられてきたマルコフ過程の推定アルゴリズムである Viterbi アルゴリズム<sup>[10],[11]</sup>がこの種の日本語の問題に有効であると考へ、特に漢字仮名変換の問題への適用を試みた。もちろん、本稿で紹介する方式が漢字仮名変換を完全に解決する訳ではなく、むしろ、この方式が従来までの単純な基本規則に加えてゆき、例外辞書や、例外コードの存在を大中に軽減する可能性があると考えた。実験もまた小規模で、結果もこのアルゴリズムで得られる精度の極値に至っている訳ではないが、この方式の有効性を示す結果が出たこと、また、この方式が漢字仮名変換以外の日本語関連の問題にも有効な手法になり得ると考へ報告する。

## 2. Viterbi アルゴリズム

Viterbi アルゴリズムとは、離散時間上の有限マルコフ過程における推定アルゴリズムであり、図1に示すようなマルコフ過程と記憶のない雑音を持つ通信路の環境において通信路の出口に

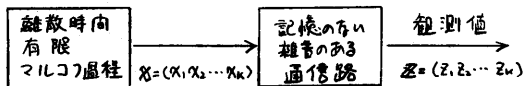
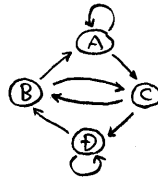


図1 一般モデル

における観測値  $z$  (以後  $z$  と略す) からマルコフ過程からの出力  $x$  (以後  $x$  と略す) を推定するものである。問題を単純化するためマルコフ過程の内部状態  $s$  と出力  $x$  を同一視すれば  $\{z_k \triangleq (x_{k+1}, x_k)\}$  通信路に記憶がないので次かいてる。

$$P(z|x) = P(z|\xi) = \prod_k P(z_k|x_k) \quad (i)$$

さて、 $z$  から  $x$  を推定するため、内部状態と時間とをそれぞれ縦軸、横軸とする有向グラフ (trellis と呼ばれる) を作ることが有効である。(図2)



内部状態

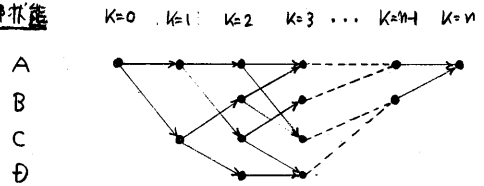


図2 trellis の例

この trellis は図1のマルコフ過程の内部状態の遷移を時間と共に表記したものであるが、重要なことは観測値  $z$  を発生させる可能性のある  $x$  がすべて trellis 上のルートとして表わされ、逆に  $z$  を発生させる可能性のない  $x$  に対応するルートは含まないように作られていることである。従って、 $z$  から  $x$  を推定することは、この trellis 上のルートを1つ選ぶことである。

$z$  を観測して、それを発生した  $x$  の最もらしいものを  $\{x^1, x^2, \dots, x^i\}$  (それぞれ trellis 上のルートに対応) から選ぶのは  $P(x^i|z)$  を最大にする  $i$  を求めればよい。ベイズの定理により、

$$P(x^i|z) = \frac{P(x^i, z)}{P(z)} = \frac{P(x^i) \cdot P(z|x^i)}{P(z)} \quad (ii)$$

であるので、結局  $P(x^i) \cdot P(z|x^i)$  を最大とする  $i$  を求めることになる。一方、(ii)式と、trellis 上の各枝が、 $x_k$  に対応しかつ  $(z_{k+1}, z_k)$  に対応している

ので、結局

$$P(x) \cdot P(z|x) = \prod_k P(z_{k+1}|z_k) \prod_k P(z_k|x_k)$$

を最大にするルートを trellis 上で求めるには、trellis の各枝に

$$-\ln P(z_{k+1}|z_k) - \ln P(z_k|x_k) \quad (iii)$$

を距離として与えれば、trellis 上の最短ルートを求めることは帰着である。そしてこの最短ルートを求めることはよく知られているように動的計画法 (ダイナミック・プログラミング) の手法を用いて容易に現実的な計算時間で解くことができる。(図3)

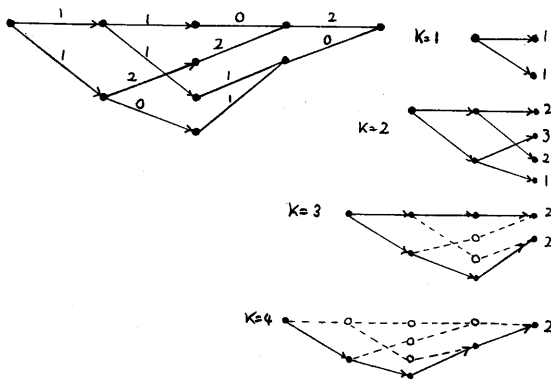


図3 動的計画法による最短ルートの探索

### 3. 漢字仮名変換の形式化

漢字仮名変換の問題をこの Viterbi アルゴリズムに形式化するには、それがマルコフ過程よりの出力文字列  $x$  が雑音のある通信路を経て文字列  $z$  に変形した図1の枠組で考える必要がある。従って、多少奇異ではあるが、そもそも単位切りされた仮名書本文  $x$  が存在しており、それが雑音のある通信路を経て漢字仮名混りのバタ書本文として観測されると考える。当然、ここでの雑音とは、 $x$  内の仮名の部分文字列が漢字文字列に確率的におき換

えられること、また、 $x$  中の単位切りの空白が失われることである。特に単位切りに関しては、応用の目的から、文節単位切り、語基レベルまでの単位切りなど様々なものも考えられる訳だが、それらのいづれもは、通信路上の雑音の設定により形式化できる。

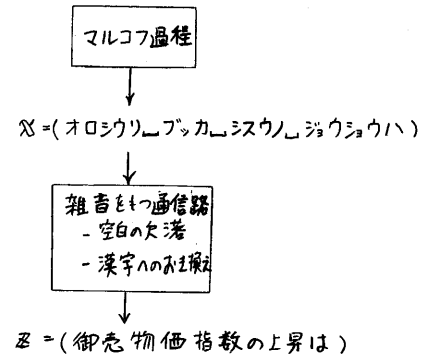


図4 漢字仮名変換のモデル

従って、この枠組での漢字仮名変換とは、観測された  $z$  を発生し得る全ての単位切りされた仮名書本文を全ルートとして含む前述の trellis を作り、各 trellis 上の枝の距離を (iii) 式により与え、その上での最短ルートを探索することになる。最短ルートが得られれば、そのルートに対応する  $x$  が、漢字仮名混りバタ書本文  $z$  の最も確からしい単位切りと仮名仮りである。

#### 3.1 雑音

図4の枠組における雑音は  $x$  上の部分仮名列が漢字、ひらがな等におき換えられるもの、また、 $x$  上の空白が欠落するものの2種を考えねばならない。前者においては、より一般的には次のような確率的雑音  $E$  を考えればよい。

$$E = \{x^* \rightarrow z^*, p(z^*|x^*)\}$$

このEは任意長の仮名列か、任意長の漢字列におき換えられるので、確率付きの国語辞典を用意することとなる。但しここでの確率とは $P(\text{漢字正書}|読み)$ である。このようなEが与えられれば観測文字列 $Z$ より、可能な $X$ をすべてのルートとして持つ trellis を作ることは容易に行える。すなわちEの逆変換 $E'$ を作れば、 $Z$ 上の部分文字列に $E'$ に含まれる変換を順次適用することですべての $X$ を含む有向グラフを作る。(図5)

$$E = \begin{pmatrix} x_1, x_2 \rightarrow z_1, z_2, & 0.3 \\ x_1, x_2 \rightarrow z_2, & 0.7 \\ x_1 \rightarrow z_1, & 1.0 \\ x_2 \rightarrow z_3, z_4, & 1.0 \\ x_3 \rightarrow z_4, & 0.9 \\ x_3 \rightarrow z_2, z_3, & 0.1 \end{pmatrix} \quad \text{有5}$$

$$Z = (z_1, z_2, z_3, z_4)$$

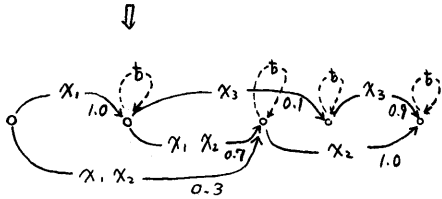


図5  $Z$ 上の $X$ を含む有向グラフ

図5において各枝の下に付した数字は雑音Eの変換の起る確率で(ii)式のオ2項に対応する。

空白が欠落する雑音は、図5に点線で示した位置に空白を発生するループを加えればよい。但しこのループは2度以上連続して回らないので、展開して、ループを含まない有向グラフに展開できる。(図6)

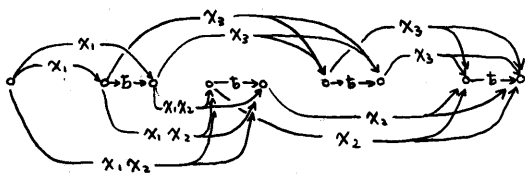


図6 空白欠落を考慮した有向グラフ

従って、図6のグラフ上の全ルートにより可能な $\{X^i\}$ 及びそれぞれに対する $P(Z|X^i)$ を得ることが出来る。

### 3.2 n-gram

雑音Eより各 $X^i$ に対する $P(Z|X^i)$ を得ることはできたが、最適な $X^i$ を選ぶためには、更に各単位切り仮名書土文の生ずる確率 $P(X^i)$ を知らねばならない。[(ii)オ] しかし、これを自然言語で分野を限らず直接的に得ることはほぼ不可能に近い。

n-gram モデルとは、文字列上の文字の発生確率がそれ以前の $n-1$ 文字に依存するとするモデルで、自然言語のふるまいを近似するのに有効であることが知られている。<sup>[12]</sup> 特に、このモデルを用いるなら、内部状態を過去 $n-1$ 文字の出力に対応させたマルコフ過程として、問題を Viterbi アルゴリズムに帰着することができる。従って、動的計画法で現実的に解くことができるようになる。

もちろん、自然言語をマルコフ過程で近似することに抵抗はあるが、文献<sup>[13]</sup>の外にも文字認識<sup>[13]</sup>、暗号解読<sup>[14]</sup>、音声認識<sup>[15]</sup>などでその有効性が実証されている。もちろん、 $m$ を大きくすれば、それと共に近似はよくなり、遠距離の文脈依存性を持つことになるが、逆に内部状態の数が $m^{n-1}$  ( $m$ はアルファベットの数)であるので計算は困難となる。また、 $P(x_n | x_1, x_2, \dots, x_{n-1})$ を計算するのための $m^{n-1}$ の数の計算が必要となるので、確率の信頼性を保つのが困難となる。

特に、漢字仮名変換などにおいては、漢字とそれの読みとの対応が遠く隔れた文脈に依存することは少ないと考えられるので小さい $n$ で十分近似が可能と思える。

n-gram で近似を行うことにより、(ii)オにおけるオ1項は $X^i = (x_0^i, x_1^i, x_2^i, \dots, x_n^i)$ に対して次のようになる。

但し  $x_0^i$  は文脈記号とする。

$$P(x^i) = P(x_0^i | x_0^0) \times \prod_{j=0}^{i-1} P(x_{j+2}^i | x_j^i, x_{j+1}^i) \quad (iv)$$

[n-gram の  $n=3$  として]

これは図6の有向グラフ上の1つのルート  $x^i$  に対応するので、図6の有効グラフより trellis を作成すれば各頂が trellis 上の枝の距離 (iii) 式の  $\alpha$  1 値に対応する。従って 3.1 の雑音で述べた (iii) 式の  $\alpha$  2 値相当の  $P(z^* | x^*)$  と合わせて、trellis 上の各枝の距離が与えられることになり、Viterbi アルゴリズムに帰着できた。

### 3.3 文脈に依存する雑音への配慮

3.1 で述べた雑音としては、仮名文字列の漢字列へのおま換え、空白の欠落の2種があったが、それらの発生確率は文脈に依存しないとしてあった。しかし、それらの発生確率が文脈に依存しないというのは、日本語の現実には合致しない。特に、文献8に見られるように単位切りのための空白と西隣の文字種の間には強い関係があるし、また、仮名列から漢字列のおま換えは特に漢字列の長土が1であるとき(おま換えが漢和辞典型のみ)西隣の文字種と強い関係がある。(例えば、訓型のおま換えを行う雑音は西隣がひらがなにおま換えられる場合だけ起り易い等)

このように雑音のおま換への発生確率がどこに近い文脈に依存して変わり得ることを無視するのは適切でない。たゞ、この雑音の依存性を表現するのは、前述の  $E$  を文脈依存型のおま換え形式、例えば  $\alpha x^* \beta \rightarrow \alpha z^* \beta, P(z^* | x^*)$ 、にあることはできない。従って、何らかの別的手段を行う必要があるが、次のように解決することができた。

前述の  $E$  はおま換え ( $x^* \rightarrow z^*; P(z^* | x^*)$ ) の集合として定義していたが、それを次のように新たに定義する。

$$E = \{ (x^*, t \rightarrow z^*; P(z^* | x^*)) \} \quad (v)$$

新たな  $E$  における  $t$  は、おま換の種類を表わす記号で、国語辞典型おま書え ( $z^*$  の長土 2 以上)、漢和辞典音型 ( $z^*$  の長土 1 で  $x^*$  が  $z^*$  の音読みである時)、漢和辞典訓型 ( $z^*$  の長土 1 で  $x^*$  が  $z^*$  の訓読みの時)、ひらがな型 ( $x^* = \text{ひ}, z^* = \text{ひ}$  等の時)、カタカナ型、特殊記号型の区別を行う。(それぞれを以後、 $W, O, K, H, N, S$  と区別する。)

雑音  $E$  を拡張すると同時に、 $n$ -gram におけるアルファベットも次のように拡張する。亦たわち、これまでの議論では  $n$ -gram は  $x_i$  亦たわち仮名文字の出現確率であったのだが、それを拡張して雑音  $E$  のおま換え規則の左辺の和集合をアルファベットとして考える。従って従来  $x_i, x_{i+1}, \dots, x_{i+n-2}$  から  $x_{i+n-1}$  の出現確率を予想するものであったが、今度は  $n$ -gram は

$$P((x_5, x_6, t_2) | (x_1, t_1), (x_2, x_3, x_4, t_2)) \quad (vi)$$

[ $n$ -gram の  $n=3$  のとき,  $t \in \{w, o, k, h, n, s\}$ ]

のように仮名の並びの推定 (上の例では  $x_1, x_2, x_3, x_4 \rightarrow x_5, x_6$ ) を行うと同時に字種の並びの推定も行う能力をも持つ。(  $t_1, t_2 \rightarrow t_3$  ) 当然空白もアルファベットの1つとして含まれるので、文字種の文脈による空白の欠落も  $n$ -gram に配慮されることになる。

この新たな雑音  $E$  と  $n$ -gram における形式化では、雑音は文脈に依存せず発生するが、それに対する補償が  $n$ -gram より完全に行われることとなる。これは図1の枠組で

通信路における文脈依存性をマルコフ過程に押しやったものであり、Viterbi アルゴリズムの枠組を用いる大々利点がある。

4. 実験

本稿で紹介した漢字仮名変換の新しい方式の有効性を確かめるために小規模の実験を行った。

4.1 確率付き漢和辞典

3.1 及び 3.2 に述べたようにこの方式では一般に雑音として国語辞典型の変換を許す訳だが、本実験では、(V)式の $Z^*$ の長さを1に制限し、確率付きの漢和辞典を用意した。(図7)

- E = ハツ, 0 → 発,  $P_1$
- ハッ, 0 → 発,  $P_2$
- パツ, 0 → 発,  $P_3$
- ハツ, K → 初,  $P_4$
- イン, 0 → 引,  $P_5$
- イン, 0 → 印,  $P_6$
- ヒ, K → 引,  $P_7$
- ヒキ, K → 引,  $P_8$
- 
- 

図7 確率付き漢和辞典

これは国立国語研究所の調査<sup>[16]</sup>を土台として若干の漢字を加え、現在約2000漢字を含んでいる。また、音便形など(発→パツ, ハッ など)に関する補正を行うと同時に、用言に用いられる漢字の送りがある(引→ひ, ひく, 取→と, とり など特に活用形)に関する補正を行っているので、通常の漢和辞典より読みの特長が多くなっている。

4.2 n-gram

一、(iii) 式の値を得るための n-gram としては、理想的には n の大きい方がよいのだが、本実験では n=3

と設定した。また実際には (vi) 式の形の 3-gram (tri-gram) を得るタータ準備の方向を軽減するため、(vi) 式の tri-gram を次式のように近似した。

$$P((X_5 X_6, t_3) | (X_1, t_1), (X_2 X_3 X_4, t_2))$$

$$\Downarrow$$

$$P(X_5 X_6 | X_1, X_2 X_3 X_4) * P(t_3 | t_1, t_2)$$

$$\Downarrow$$

$$P(X_6 | X_4 X_5) * P(X_5 | X_3 X_4) * P(X_4 | X_2 X_3) * P(X_3 | X_1 X_2) * P(t_3 | t_1, t_2) \quad (vii)$$

従って、(vi) 式を得るために、字種の tri-gram と仮名文字の tri-gram が必ずであった。

字種の tri-gram は信頼のおけるタータが手に入らなかつたので、文献8の坂本氏の調査を基にして簡単なモデルによる推定値を加え代用した。

一、仮名文字の tri-gram は、仮名書きして1200~1500字程度の文章を36ヶ集め、単位切り仮名書きし、仮名文字の3つ組を集計して作成した。文章の分野としては、新聞の社説約半分、コンピュータ関連のマニュアル、雑誌記事が約1/4、残りが社内規則である。ここでの単位切りは文節切りに加えて、4文字以上の複合語内を2文字トしくは3文字の語基に単位切りを行うことを行った。(図8) 従って後に示す結果の単位切りにもその方式が反映している。

チホウ ファンケン ノ リネワ イカ  
 タ イ 17 ジ 子ホウ イト 子ヨウサリイカ、アタラシ  
 プリカダニツイテ トウジツアツコ マトメタ。9カ ツチュウニ  
 トウジツアツコ タイヨウハ、クニ、子ホウヲ ツクシタ  
 オキテ ハンラトツツテイル ソノタメトシ タイカク ノチ  
 キ ヨウエイ カクノヲ 子ホウニ イシヨウ、コソコ 補  
 子ホウ テ サキ カルヲ ハイシ トウゴウスル サラニ 子  
 キ ヨウエイ タイイ キ ヨウエイ イソカイ イト キ  
 コホフニキ コホイニ テイキサレイル

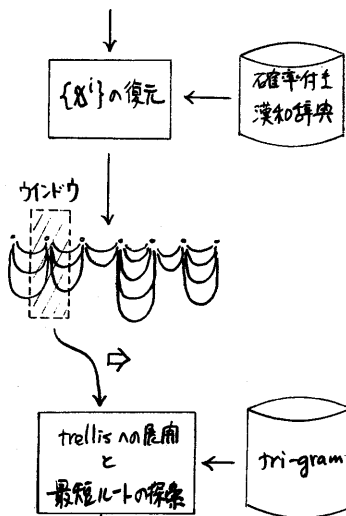
図8 訓練用タータ

特に単位切り工ねた訓練用データから仮名の3つ組を集計する際に、幼音、長音、濁音、半濁音などは先行する文字が限定されるので、独立した仮名文字とは扱われなかった。従って、むしろ、もう単位の tri-gram でありといえる。

### 4.3 解析

解析プログラムは PL/I で約 3000 ステップで図9の構造である。入力に対して、まず図7の確率付生漢和辞典を引き、図6の形の{x}を含む有向グラフを作成する。次に有向グラフの図9のウィンドウ部分を trellis に部分的に展開する。(trellis 全体を展開するのはメモリーを食うため) このウィンドウを有向グラフの前から後ろへ移動すると同時に動的計画法により最短ルートが得られる。図10に、入力列「再上昇に転じな米金利」に対して得られる図6形式の有向グラフを例として示す。

入力列: 再



出力列: 再  
図9 解析の流れ

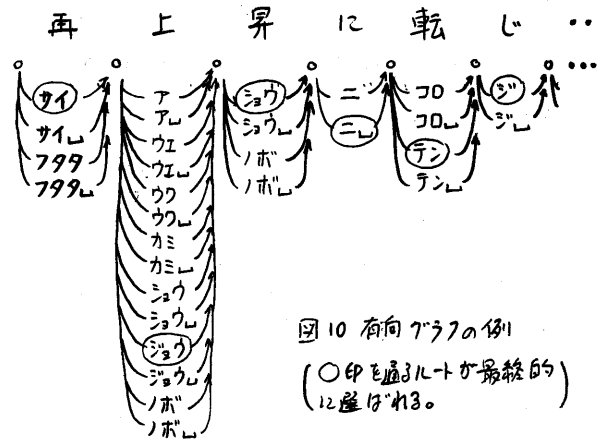


図10 有向グラフの例  
(○印を辿ると最終的に選ばれる。)

### 5 結果と考察

新聞の社説を適当に選んで本サイトで漢字仮名変換を試み、例の一部を図11に示す。

(図11)

引き続き高水準を続けている。一月の卸売物価指数は完成品・季節調整済

( ) ( ) 引(ヒ)き(キ)続(ツツ)き(キ) 高(コフ)水(スイ)準(ジュン)を(ヲ) 続(ツツ)け(ケ)て(テ)い(イ)る(ル)。 ( ) - (イ)チ)月(カ)ツ)の(ノ) 卸(オロシ)売(ウリ) 物(フツ)価(カ) 増(シ)強(スツ)は(ハ) 売(カン)成(セイ)品(ヒン)。 ( ) 準(キ)節(セツ) 調(チョウ)整(セイ)済(ザイ)。 ( )

ベースで年率約20%も上昇した。新方式の採用後、しばらく

( ) ( ) ベー(ペー)ス(ス)で(テ) 年(ネン)率(リツ)約(ヤク) D( D) 。 ( ) も(モ) 上(ジョウ)昇(ジョウ)し(シ)を(タ)。 ( ) 新(シン)方(ホウ)式(シキ)の(ノ) 採(サイ)用(ヨウ)後(ゴ)。 ( ) し(シ)ば(ハ)ら(ラ)く(ク)。 ( )

落ち替っていた通貨供給量や銀行貸出しの増勢にも、昨年末から

( ) ( ) 落(オ)ち(チ)替(ツ)い(イ)て(テ)い(イ)を(タ) 通(ツウ)貨(カ) 供(キョウ)給(キョウ)量(リョウ)や(ヤ) 額(ガン)行(コウ)貨(カ)し(シ)出(ダ)し(シ)の(ノ) 増(ゾウ)勢(セイ)に(ニ)も(モ)。( ) 昨(サク)年(ネン)末(マツ)か(カ)ら(ラ)。 ( )

再び強まる傾向がみられる。FRBとしては一月の卸売物価の動きが

( ) ( ) 再(ワタ)び(ヒ) 強(ツヨ)ま(マ)る(ル) 傾(ケイ)向(コウ)が(ガ) み(ミ)ら(ラ)れ(レ)る(ル)。 ( ) A(A)と(T)し(シ)て(テ)は(ハ) - (イ)チ)月(カ)ツ)の(ノ) 卸(オロシ)売(ウリ) 物(フツ)価(カ)の(ノ) 動(ウゴ)き(キ)が(ガ)。 ( )

確認した物価の悪化持続、原油の値上がりそれがさらに強めそうな

( ) ( ) 確(カク)認(ニン)し(シ)を(タ) 物(フツ)価(カ)の(ノ) 悪(トウ)化(セイ) 持(シ)続(ジツ)。 ( ) 原(ガン)油(ユ)の(ノ) 値(ネ)上(ア)が(ガ)り(リ)が(ガ) そ(ソ)れ(レ)を(ラ)さ(サ)ら(ラ)に(ニ) 強(ツヨ)め(メ)そ(ソ)う(ウ)を(タ)。 ( )

情勢を前に、改めてここで政策方針を明示する必要を感じたのであろう。

( ) ( ) 傾(ジョウ)勢(セイ)を(ラ) 前(マエ)に(ニ)。 ( ) 派(アラタ)め(メ)て(テ) こ(コ)こ(コ)で(テ) 強(セイ)化(クワ)方(ホウ)針(シン)を(ラ) 明(メイ)示(シ)す(ス)る(ル) 点(ヒツ)聖(ヨウ)を(ラ) 確(カン)じ(ジ)を(タ)の(ノ)で(テ)あ(ア)る(ロ)う(ウ)。 ( )

アフガニスタン問題の発生に伴う米国の軍事生産の拡大、朝鮮半島の

( ) ( ) ア(ア)フ(フ)ガ(ガ)ニ(ニ)ス(ス)タ(タ)ン(ン)領(モン)土(ダ)の(ノ) 売(ハツ)生(セイ)に(ニ) 件(トモナ)う(ウ)米(ペイ)国(コク)の(ノ) 軍(ガン)費(ジョ) 生(セイ)産(サン)の(ノ) 増(カク)大(ダ)イ。 ( ) 戦(セン)時(リョク) 物(フツ)価(クワ)の(ノ)。( )

表1にその社説全体における結果を示す。

a	総文字数	1072
b	総漢字数 (漢数字は「一」以外は含まない)	461
c	誤って読まれた漢字数	23
d	読みぶり正答率(c/b)	95.01%
e	単位切り位置数	273
f	余分に挿入された空白数	7
g	単位切り忘れの個数	16
h	単位切り正答率 $(\frac{e-f-g}{e})$	91.58%

(注) ここの正しい単位切り位置とは文節内を切り、複合語内語基内を切りである。

表1 結果

表1の結果は、訓練データの少ない、(約45K仮名文字)に対して十二分に満足できるものと考えられる。訓練データを蓄積することにより、まだ精度は向上すると思われる。

逆に、この実験もしくは形式化の限界を示す誤りも若干見られる。例えば「国家的計画、」が「ツッカケキケイガ、」のように変換された。これは(vi)のtri-gramを(vii)のように近似したことに起因していると考えられる。あるいは(vii)式の近似によるtri-gramで評価される仮名列がすでにそれらの変換されるべき字種を忘れてしまっている。従って、この例でも「計画」より得られる仮名列「ケイカク」と「ケイガ」の比較で、後者が「系が」、「計が」、などと互換的に頻度の高い別の同音表現のために優先してしまう、ている。

別の例で興味深いのは「歩合」かいつでも「フコウ」と呼まれてしまうことである。これも「歩合」から得られる「フパイ」、「フコウ」の比較において、「符号」、

「符号」、「巨号」、「富豪」などのより頻度の高い同音語にひきかかれていると思われる。

単位切りが91.58%とあまりよくないのは主として長い複合語分解の誤りによる。例えば、「連邦準備理事會」は「連邦準備理事會」のように分割された。これは、現在のモデルの複合語分解が、漢語は2文字単位が多いなどの複合語に固有のいくつかの常識を用いず、単に空白を仮名列内、字種内に含んだ場合のtri-gramによる評価だけで行っているからと考えられる。これは近い将来の課題である。

おわりに

本稿で紹介したViterbiアルゴリズムによる漢字仮名変換は実用的な観点からは精度が十分でないとの見方もある。しかし、まな書まにも述べたように、この上に例外辞書、例外コードなどを積み上げるための土台として十分有効と考える。また、モデル自体にもまだ幾多の本質的改良の可能性がある。実験も小規模でこのモデルの性能を十二分に引き出していない点もある。しかし、本稿で紹介したViterbiアルゴリズム及びそれでの漢字仮名変換の形式化は全く新しく、他の関連問題でも十分に有効な手段であると考えたので報告を行った。

なお、最短ルートの探索パスのプログラム作成で、当センターの客員研究員瀧川清氏(早稲田大学)に協力いただいた事をご感謝する。

参考文献

[1] 藤崎他、日本語文書処理システム「ことばま」-概念と設計思想、東京カインテックセンターレポート、N:9318-1512、1980  
 [2] 藤崎他、日本語文書処理システム「ことばま」の仮名漢字変換、情報処理論文誌、vol.23, No.1, 1982 (予定)



- [3] 藤崎、諸橋、ことば文書処理システムの文書検索機能, IBMレビュー 85, 1981 (予定)
- [4] 大河内他、仮名漢字変換のための文法解析、情報処理計算言語学研究会資料 25-4, 1981
- [5] 田中、漢字がなまじり文を全文カナ書き、D-マ字書きに変換するシステムについて、国立国語研究所報告(34), 1969
- [6] 中野、言語研究のための日本語テータ入力システム、日本語情報処理シンポジウム報告集、1978
- [7] 荒木他、JICSTの実用的全自動漢字-カナ変換システム、K-KACS 12について、情報処理、Vol 20, No. 10, 1979
- [8] 坂本、文節の認定、日本語情報処理シンポジウム報告集、1978
- [9] 長尾、他、国語辞書の記憶と日本語文の自動分割、情報処理、Vol 19, No. 6, 1978
- [10] G. D. FORNEY, JR., The Viterbi Algorithm, Proceedings of IEEE, Vol. 61, No. 3, 1973
- [11] A. J. Viterbi, Error bounds for convolutional codes and an asymptotically optimum decoding algorithm, IEEE Trans. Information Theory, vol. IT-13, 1967
- [12] C. E. Shannon, Prediction and Entropy of Printed English, Bell Sys. Tech. J., Vol. 30, 1951
- [13] J. Raviv, Decision Making in Markov Chains Applied to the Problem of Pattern Recognition, IEEE Trans. Information Theory, Vol IT-3, No. 4, 1967
- [14] L. R. Bahl, et. al., Disambiguation of Stenotype and Other forms of Ambiguous Text by Use of Contextual Statistics, personal memo, IBM T.J. Watson Res. Ctr, 1980
- [15] J. Jelinek, Continuous Speech Recognition by Statistical Methods, Proc. of the IEEE, Vol. 64, No. 4, 1976
- [16] 現代新聞の漢字、国立国語研究所、秀英出版、1976