

# 単語分割と単語選択に関するかな漢字変換法の評価に関する一考察

荒木哲郎 倉野真樹 山田和義 古川貴康 小越康宏 大場達哉

福井大学工学部知能システム工学科

{ araki , kurano , ogoshi } @human.his.fukui-u.ac.jp

## 1. はじめに

従来、べた書きかな文のかな漢字変換法としては、文節単位に分割する方法には、2文節最長一致法<sup>[1]</sup>、文節数最小法<sup>[1]</sup>、かな文字のマルコフ連鎖モデルを用いた方法<sup>[2]</sup>などが提案されており、また、単語選択には連語解析を用いた方法<sup>[1]</sup>、格フレームを用いる方法<sup>[1]</sup>、漢字かな混じり文字のマルコフ連鎖モデルを用いる方法<sup>[3]</sup>、確率モデルを用いる方法<sup>[4]</sup>、読み情報を用いる方法<sup>[5]</sup>などがある。日本文音声入力のかな漢字変換を考えると、さらに変換精度の向上が望まれる。本論文では、仮単語境界推定を用いたべた書きかな文のかな漢字変換法を提案するとともに、かな漢字変換法を単語分割と単語選択の観点から方式の異なる2つの方式について、新聞記事データを用いた実験を行って、かな漢字変換精度の評価と辞書アクセス回数の比較を定量的に行う。

## 2. かな漢字変換法のモデル

### 2.1 べた書きかな文の単語分割とかな漢字変換候補の絞り込み処理を並行して行うかな漢字変換法<sup>[3]</sup>

2.2で提案されるかな漢字変換法と比較する上で、ここでは<sup>[3]</sup>で示されているかな漢字変換法を方法Aとし、手順を以下に示す。

【ステップ1】べた書きかな文節に対し、単語辞書の読み見出しにより検索し、先頭から順に部分列に一致する漢字かな単語候補を全て抽出する。

【ステップ2】ステップ1で得られた各部分列に対する漢字かな単語候補を相互に結合して構成

される漢字かな文節候補を全て生成する(漢字かな文節ラテイスと呼ぶ)。

【ステップ3】ステップ2の各漢字かな文節候補に対して、漢字かな文字の2重マルコフ連鎖モデルを用いて評価し、最大の連鎖確率値を持つ単語候補列を最尤な単語候補列として選択する。

この手順に従って得られる漢字かな文節候補ラテイスの例を図1に示す。

### 2.2 仮単語境界を用いたかな漢字変換法

べた書きかな文の文節単位の分割として、仮文節境界を推定する方法<sup>[2]</sup>が提案されていたが、ここでは、かな表記の単語マルコフ連鎖モデルを用いて仮単語境界を推定する方法を提案する。この方法を方法Bとし、以下に手順を示す。

【ステップ1】べた書きかな文節に対して、同じ読みを持つ複数の単語候補を一つかな表記で表し、文節を構成するあらゆるかな表記の単語候補列の組み合わせを網羅的に全て抽出する(かな単語候補ラテイスと呼ぶ)。

【ステップ2】ステップ1で抽出されたかな各単語境界候補列に対して、かな表記の2重単語マルコフ連鎖モデルを用いて評価し、最尤な単語境界候補を決定する。

【ステップ3】ステップ2で求めたかな単語境界に対して、単語辞書を用いて、同じ読みを持つ漢字表記の単語候補を全て抽出する。

【ステップ4】漢字かな文字のマルコフ連鎖モデルを用いて最尤な漢字かな変換候補を決定する。

仮単語境界推定を用いたかな単語候補ラティスの例を図2に示す。

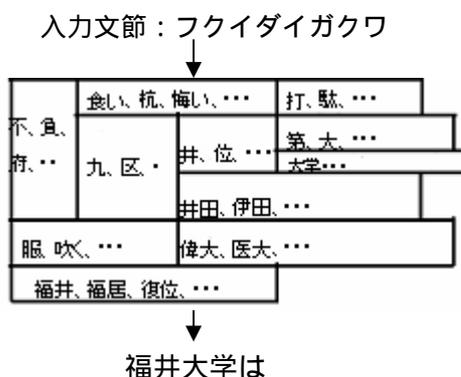


図1 方法Aによる漢字かな単語候補ラティスの例

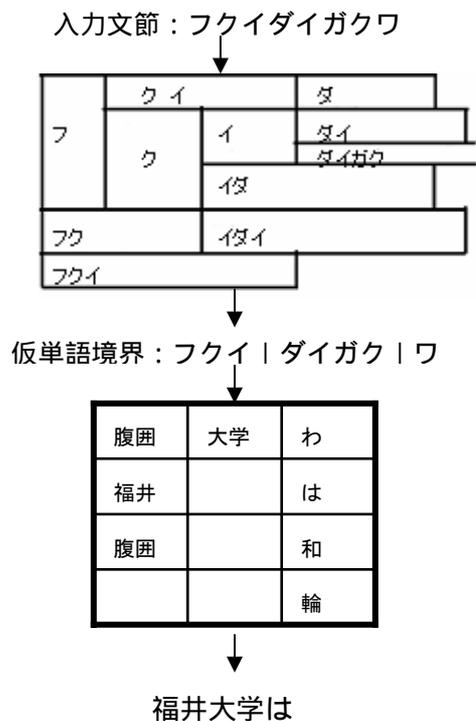


図2 方法Bによるかな単語候補ラティスと漢字かな単語候補ラティスの例

## 2.3 マルコフ連鎖モデルの違いによるかな漢字変換精度の評価モデル

マルコフ連鎖モデルの違いによるかな漢字変換方法を次の観点から分類し、表1のような評価モデルを設定する。

べた書きかな文字列を、単語辞書引きを用い

て単語単位に分割する際に、どのような情報を用いて分割するか。

で決定された仮名文字列の単語分割に対し、単語辞書を用いて単語候補を選択する際に、どのような情報を用いて選択するか。

このような2つの処理に対して、次のようなマルコフ連鎖モデルを用いてモデルを設定する。

情報の単位が文字(C)か、単語(W)か  
マルコフ連鎖の次数が、1重(1)か、2重(2)か  
字種が、かな文字表記(KN)か、漢字かな混じり表記(KK)か

例えば、2重文字単位のかな文字連鎖モデルを2KNCモデルと表すことにする。

表1 かな漢字変換方式の評価モデル

モデル	単語分割	単語選択
A-1	2 KKC	2 KKC
A-2	1 KKW	1 KKW
A-3	2 KKW	2 KKW
B-1	2 KNW	2 KKC
B-2	2 KNW	1 KKW
B-3	2 KNW	2 KKW

## 3. 実験結果

### 3.1 実験条件

#### (1) 入力データ

文の種類：日本経済新聞記事

字種：べた書きかな文節

総文節数：標本外データ 668 文節(100 文)

文節平均 6.9 文字，文平均 47.4 文字

#### (2) 使用辞書

40 万語の単語辞書

単語境界推定に用いる文節単位のかな表記単語 2 重マルコフ連鎖確率辞書

かな漢字変換候補の絞り込みに用いる漢字かな混じり文字の 2 重マルコフ連鎖確率辞書

### 3.2 実験結果と考察

#### (1) 単語分割推定精度の比較

方法 A と B による単語分割の精度を、表 1 に示されるマルコフ連鎖モデルを用いて絞り込まれたときの第 1 位候補方式について比較したものを表 2 に示す。同表から、方法 B の方が高いことがわかる。これは単語分割に対して、かな表記の単語候補の数が、漢字かな混じり表記の単語候補の数より極めて少ないことによる効果が現れていると考えられる。

また、方法 B の場合についてかな文字表記の単語分割精度の 10 位までの累積正解率を図 3 に示すと、10 位までの累積正解率は 99.8% と非常に高い値が得られことがわかった。

表 2 第一位正解候補の単語分割精度

A - 1	62.5%
A - 2	82.9%
A - 3	83.6%
B	98.9%

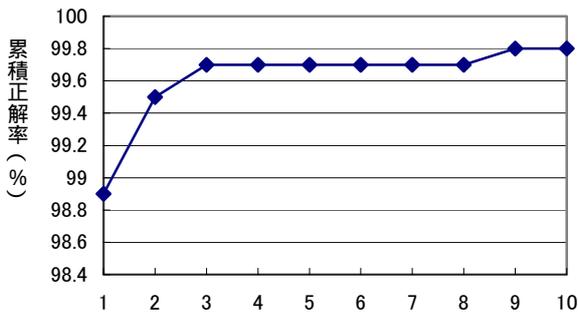


図 3 方法 B のかなの単語境界推定精度

#### (2) かな漢字変換精度

方法 A と方法 B による、べた書きかな文節のかな漢字変換結果を図 4 に示す。同図より、第一位正解率で見ると、方法 B-2 が一番高く、累積 10 位までの累積正解率では方法 B-3 の精度が一番高い結果となった。いずれも、方法 A と比べ

て、方法 B の方が高い精度がえられることが分かった。これは、方法 B で高い単語分割精度が得られるため、単語選択で方法 B の方が、方法 A より選択候補数が少ないことが原因と考えられる。

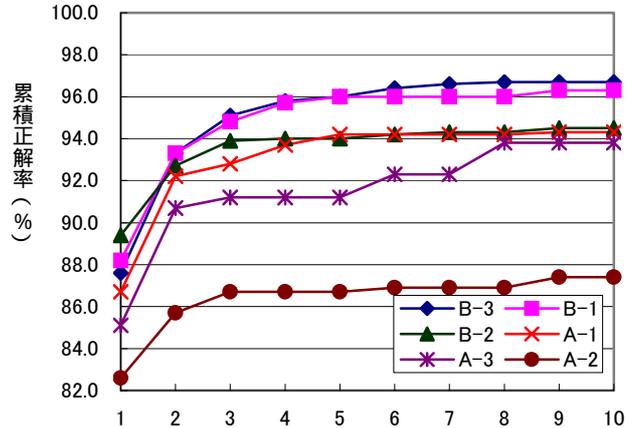


図 4 かな漢字変換モデルの変換精度の比較

#### (3) かな漢字変換候補の第一位文字の適合率と再現率

ここでは、かな漢字変換モデルにおける最尤候補と正解候補との最長共通部分列 (longest common subsequences) [6] の文字数に基づく再現率と適合率を表 3 に示す。同表より、最長共通部分列の評価では、方法 B の方が適合率及び再現率共に若干高い値を示したが、ほとんど大差は見られなかった。

表 3 最長共通部分列による各種変換モデルの最尤候補と正解候補の適合率と再現率

モデル	適合率(%)	再現率(%)
A-1	94.2	93.4
A-2	94.2	93.4
A-3	94.4	93.4
B-1	94.9	95.2
B-2	95.5	95.4
B-3	94.5	94.6

#### (4) 漢字変換処理に要する辞書アクセス回数の比較

処理時間を評価する上で、方式 A と B でかな漢字変換に必要な単語候補の組み合わせ数の比較結果を、1重単語マルコフ連鎖モデル(1KKW)を用いた場合について図5に示す。

また、2重文字マルコフ連鎖モデル(2KKC)を用いた場合の辞書アクセス回数の比較を図6に示す。同図より、方法 B がかな文字列の長さが7の場合で、約100倍アクセス回数が少ない結果となり、処理時間が大幅に削減されることがわかった。

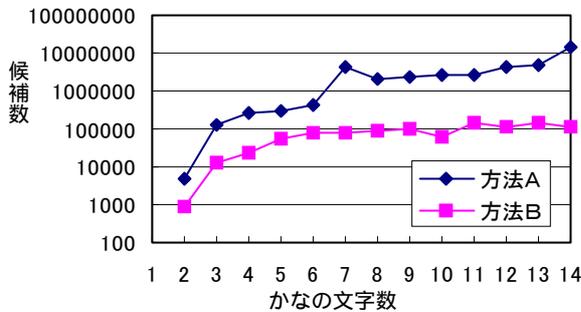


図5 1重単語マルコフ連鎖モデルの単語候補の組み合わせ数の比較

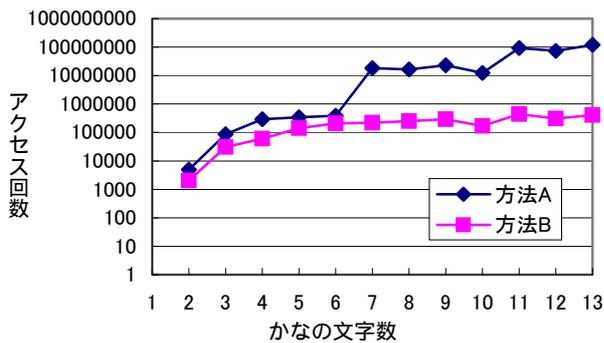


図6 2重マルコフ連鎖モデルを用いた場合の辞書アクセス回数の比較

## 5. まとめ

かな漢字変換法において、単語分割法と単語選択法の異なる2つの方法に対し、(1)情報単位が文字単語か、また(2)マルコフ連鎖モデルの次数

が1重か2重か、及び(3)字種がかな文字表記か漢字かな混じりかについて組み合わせ、6通りのモデルを設定した。実験により、文節単位及び最長共通部分列による比較では両方ともに余り大差は見られなかったが、仮単語境界を用いた後で単語選択を行う方法が、単語分割と単語選択を同時に行う方法に比べて若干高いことがわかった。

また、辞書へのアクセス回数は、7文字以上のとき約 $1/10^2$ に縮小されることがわかった。

このことから、単語分割を行った後で単語選択を行う方法が変換精度及び処理時間の上から優れていることがわかった。

## 文献

- [1] 齊藤裕美, 川田勉: “ 仮名漢字変換アルゴリズム ”, 進学誌, vol. 70, No. 8, pp679-687, (1986)
- [2] 荒木哲郎, 池原悟, 橋本昌東, 三品尚登: “ 2重, 3重のマルコフ連鎖モデルを2段階に使用したべた書き仮名文の文節境界推定法 ”, 信学論, Vol. J83-D-II, No. 12pp2745-2754, (2000.12)
- [3] 村上仁一, 荒木哲郎, 池原悟: “ 日本文音節入力に対して2重マルコフ連鎖モデルを用いた漢字仮名交じり文節候補の抽出精度 ”, 信学論, Vol. J75-D-II, No. 1, pp11-20, (1992.1)
- [4] 森信介, 土屋雅稔, 山地治, 長尾真: “ 確率的モデルによる仮名漢字変換 ”, 情処論, Vol. 40, No. 7, pp2946-2953, (1999.7)
- [5] 荒木哲郎, 池原悟, 真田陽一, 横川秀人: “ 読み情報を用いた仮名漢字変換の精度向上効果の推定 ”, 信学論, Vol. J84-D-II, No. 2, pp351-361, (2001.2)
- [6] Aho. A. V: “ 文字列中のパターン照合のためのアルゴリズム ”, コンピュータ基礎理論ハンドブック, Vol. 1, 形式的モデルと意味論, pp. 263-304, Elsevier Science Publishers (1990)