

5 L-8

仮名漢字変換における 変換手法と変換精度についての比較実験

酒井貴子, 下村秀樹, 並木美太郎, 中川正樹, 高橋延匡
(東京農工大学 工学部 電子情報工学科)

1.はじめに

計算機に対する入力手段として仮名漢字変換が定着したことから、ワードプロセッサなどを使用して文章を作成するユーザが増加した。そこで、(1)仮名漢字変換手法の研究[1]、(2)学習機能や意味、文法情報の利用の研究[2]などが行われ、変換精度の向上が図られてきた。より良い変換を実現するためには、これらに対する評価を行う必要があるが、変換手法や各機能の性能を同一システム上で評価した例は少ないという現状にある。

そこで、我々はこうした仮名漢字変換手法に対する評価が容易に行える、実験環境としての仮名漢字変換システム（以下、本システムとする）を開発した[3]。我々はこれを用いて、(1)変換手法間での性能の差、(2)変換手法と機能との組合せによる効果の度合い、などを同一条件で定量的に評価して、問題点を明らかにしたいと考えた。本稿では、本システムを用いて行った、複数の変換手法と変換精度の比較実験について述べる。

2. 最尤候補選択法による仮名漢字変換システム

仮名漢字変換では、複数通りの変換結果が考えられるため、複数の候補から一つを選択するための評価基準を設ける。変換手法はこの評価基準に相当するが、手法には最長一致法のように「文頭の単語の長さ」という局所的な評価を行う場合と、文節数最小法のように変換結果全体が評価対象になり、すべての変換結果を作成してみないと最尤候補が得られない場合がある[1]。本システムは、様々な変換手法や機能の性能評価を行うことを目的としているため、どのような手法でも実現可能にする必要がある。そこで本システムでは、次の2点を特徴とした。

- (1) 初回の変換で考えられるすべての候補を作成する。
- (2) 変換結果を作成する処理と変換結果を選択する

処理を完全に分離して実現する。

(1)は、変換時間やメモリ容量など物理的な制約が予想されるが、実験によって実現できる範囲で収まることが確かめられた[3]。また(2)により、変換結果を選択する処理を作成して替える手間だけで、異なる変換手法を実現できる。

本システムは、変換結果の尤度を評価値で表現し、すべての変換結果に対して評価値を付け、評価値の高い順番に変換結果として選択する方法（最尤候補選択法）に

基づく[3]。この場合の評価値は、変換結果を構成する各単語と単語間の接続関係に対して付くものとした。したがって、変換手法は変換結果に評価値を付ける処理として実現される。こうした方式に基づき、変換手法の性能評価を行いややすい実験環境としての仮名漢字変換システムを実現した。

3. 仮名漢字変換の評価用ツールキット

仮名漢字変換手法や各機能に対して定量的に評価を行うためには、大量の文章を仮名漢字変換して変換精度を測定する必要がある。今回は、変換精度を次式の変換率(%)によって示すこととした。

$$\text{変換率} = \frac{\text{期待する結果が1回の変換で得られた数}}{\text{入力データ総数}}$$

変換率の測定を人手で行うには限界がある。そこで我々は、漢字仮名混じり文を仮名文に変換するツールを作成し、文書ファイルから半自動的に入力用の仮名データを作ることができるようにした。

また、初回の変換ですべての変換結果を作成するという本システムの特徴を利用して、変換結果の中から正しい変換結果を探索するツールを作成した。このツールは、正しい変換結果が探索されるまでの変換回数、文節移動回数などが計測でき、自動的に変換率を測定することができる。これらの仮名漢字変換の評価用ツールを図1のように使用することで、大量のデータに対する実験を可能にした。

4. 仮名漢字変換手法と変換精度の比較実験

こうした実験環境で、複数の変換手法について変換精度を測定する実験を行った。詳細を次に示す。

実験用ベンチマークテキスト

- (1) 情報処理学会論文誌から2報[4][5]
- (2) 情報系の卒業論文1報

(1)はB5判で2段組約7~8ページ、(2)は約1400文字/ページの本文64ページを使用し、(1),(2)は別々に実験を行った。また、実験時の入力単位は句読点までの長さを1データとした。これは、句読点を挟む文字間では変換結果に曖昧さが生じないことから、最長の変換単位に相当すると考えたことによる。

Measurements of the accuracy of various translation algorithms for Kana to Kanji translation

Takako SAKAI, Hideki SHIMOMURA, Mitarou NAMIKI, Masaki NAKAGAWA and Nobumasa TAKAHASHI

Department of Computer Science, Faculty of Technology, Tokyo University of Agriculture and Technology

実験対象の変換手法

今回は、次に示す三つの代表的な仮名漢字変換手法を実現し、実験を行った。()内は略称とする。

- (1) 最長一致法(最長)
- (2) 2文節最長一致法(2文節)
- (3) 文節数最小法(文節数)

また、これらに学習機能を組み合わせたときの効果を測定するために、次の二つの学習方法を実現した。“学習”とは、「それまでの変換結果がそれ以降の変換結果に影響を与えること」と定義する。

- (1) 最近使用語優先学習(学1)

ユーザによって使用された語(学習語)を記憶し、最近使用された語を含む変換結果から選択する。

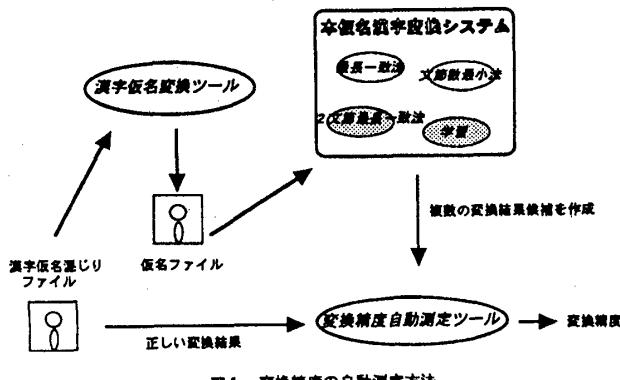
- (2) 学習語数優先学習(学2)

学習語が最も多く含まれる変換結果を選択する。

上述した三つの基本手法と二つの学習方法を組み合わせ、さらに手法と学習の優先関係を入れ換えることで、10種類の評価関数を実現した。これらに対して、基本手法単体からの変換率の伸びを調べた。

実験方法

3章で述べた仮名漢字変換評価用ツールを図1のように組み合わせ、変換精度の測定を自動化した。



実験結果

上述した各評価関数に対してベンチマークテキスト(1), (2)を独立に実験し、測定した変換率の平均値を図2に示す。図2のグラフは白い部分が各基本手法単体で実験を行ったときの変換率、黒い部分が各評価関数の学習による変換率の伸び率を示す。次に、これらの結果を踏まえて考察を行う。

① 三つの基本手法単体での変換率の差は、最大でも7.7%と以外に小さい。ただし、学習を行った場合の変換率を比較すると、差は29.9%(図2 最長+学1と文節数+学1)にも及ぶ。これは、基本手法だけでは文節の区切り方しか限定できず、同音語の選択能力がないため、効果があまり現れないことを数値的に示した結果といえる。

なお、手法なしでの変換率は9.3%であり、基本手法

を使用することによって9~17%の変換率が向上するという結果も得られた。

② 学習を行うと、最大で43.7%もの変換率が向上した。また学1と学2の学習方式では、学1の方式が変換精度の向上効果が大きかった。これは、ユーザが最近使用したという情報が有効であることを示した結果といえる。

③ 手法と学習の優先関係については、学習よりも手法を優先した場合が圧倒的に変換率が高い結果となった(例 文節数+学1は69.7%に対して、学1+文節数は33.7%と1/2以下)。また、手法によって文節の区切り方が決定された上で、学習情報を利用しないと、学習が悪影響を及ぼす(例 学2+2文節、学2+文節数ともに変換率低下)可能性があることも明らかになった。

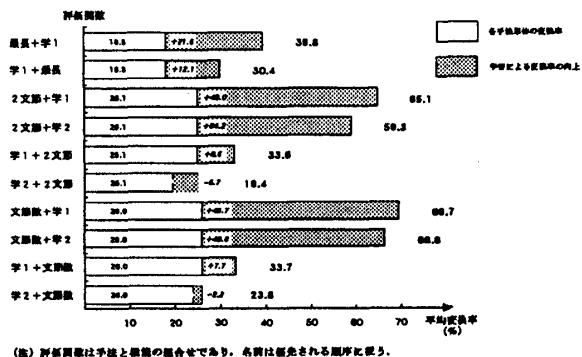


図2 各変換手法における変換率

5. おわりに

本稿では、仮名漢字変換の基本手法と学習について、変換率に与える効果を測定した実験について述べた。今後は、(1)文章量を増やすこと、(2)異分野の文章について実験を行うこと、が課題として残されており、こうした定量的な評価を行うことによって、仮名漢字変換技術の問題点の解決を図りたいと考えている。

参考文献

- [1] 吉村他：最長一致法と文節数最小法について、情処人工知能と対話研究会報告 24-1, 1982
- [2] 酒井他：日本語ワードプロセッサの仮名漢字変換における変換処理と精度についての考察、情処 HI 研究会資料 35-10, 1991
- [3] 酒井他：仮名漢字変換における最尤候補選択アルゴリズムの実験、情処第44回全国大会論文集4P-12, pp.191-192, 1992
- [4] 下村他：最小コストパス探索モデルの形態素解析に基づく日本語誤り検出の一方式、情処論文誌, Vol.33, No.4, pp.457-464, 1992
- [5] 下村他：人間の誤り検出能力と誤り検出機能の効果に関する実験、情処論文誌, Vol.33, No.12, pp.1607-1617, 1992