

E-019

ニューラルネットワークを用いた携帯端末向け日本語入力手法の有効性について

Effectiveness for Fast Japanese Input Method Using Neural Network on Mobile Terminal

鈴木 悟史[†] 松原 雅文[†] Goutam Chakraborty[†] 馬淵 浩司[†]
 Satoshi Suzuki Masafumi Matsuhara Goutam Chakraborty Hiroshi Mabuchi

1. はじめに

現在、日本において携帯電話は、電子メールや様々なアプリケーションを利用することが出来る端末として発達しており、携帯電話上での日本語文入力の機会と必要性が増大している。しかしながら、携帯端末はそれ自身の特性から小型であることが求められ、通常のフルキーボードより遥かに少ないキー数となることは必至である。このため現在の携帯電話における日本語入力方式は、キー数の少なさを補うために多くの打鍵数を必要とし、迅速な入力が困難であるという問題を抱えている。

この問題を解決するために、「ニューラルネットワークを用いた携帯端末向け日本語入力手法」[1]が提案されている。この提案手法では、文字情報縮退方式という入力方式を利用し、変換の際に問題となる曖昧性をニューラルネットワークで解決しようとした。この提案手法の実験結果は、有効性を示唆するものであったが、入力が単語単位という制約や、ニューラルネットワークの規模が定まらないという問題があり、実用化は困難であると考えられる。

そこで、本研究では入力を固定長で分割した数字列とすることで、ニューラルネットを利用しながらも、先行研究の問題点を解消し、より実用に適した、携帯端末向け日本語変換手法を提案する。本稿では、オープンデータによる実験を行った結果から、本手法の有効性について述べる。

2. 先行研究

2.1 文字情報縮退方式

文字情報縮退方式は図1のようなキー配置であり、これは現在主流である文字循環指定方式と同じである。このため、利用者は新たにキー配置を覚え直す必要が無い。

具体的な入力方法を「大会(たいかい)」という例を使い説明する。「たいかい」を文字情報縮退方式により入力する場合、4121の順にボタンを打鍵する。そして入力された数字列を「大会」に変換する。このように、入力したい仮名が含まれる数字を1回、又は濁点半濁点の場合は2回打鍵するだけで、1文字の入力が可能である。よって従来の文字入力方式に比べ迅速な入力が可能となる。

しかし、母音情報が縮退されているため入力数字列が非常に多くの曖昧さを含むという問題がある。「たいかい」を表わす4121は、他に「とうけい」や「つうこう」など $5^4 = 625$ 種類の仮名文字列に対応している。また、日本語は漢字への変換も必要となるためさらに多くの量の候補が存在する。このため、縮退された数字列を日本語に変換する手法(以下、数字漢字変換手法)[2]が、文字情報縮退方式を用いる上で重要となる。

1 "あ行+あ行" "ー"	2 "か行"	3 "さ行"
4 "た行" "っ"	5 "な行"	6 "は行"
7 "ま行"	8 "や行" "やゆよ"	9 "ら行"
* "濁点" "半濁点"	0 "わ行"	# "句読点"

図1: 文字情報縮退方式のキー配置

2.2 ニューラルネットワークを用いた変換手法

曖昧さを解決するために、先行研究ではニューラルネットワークを利用し、その高度な学習機能を活用することで数字漢字変換を行った。誤差逆伝播法のニューラルネットワークで学習を行い、入力に単語の数字列と、その数字列の前後に入力された数字の出現頻度を与えることで、日本語の文字コードが出力される仕組みである。

実験はクローズドデータで行われ、その結果、ノード単位では97.3[%]、単語単位では62.0[%]の正解率となり、有効性が示唆されたが、入力される数字列が単語単位であるという前提のもとに実験が行われていた。しかし、実際の入力においては必ずしも単語単位で区切って入力を行うとは限らない。よって、どのような長さ、区切りの数字列にも対応することが可能な変換手法が必要である。

3. 提案手法

本手法における数字漢字変換の流れを図2に示し、図3を例に挙げて説明する。まず、はじめに入力数字列を受け取り、固定長による分割処理を行う。分割処理は、はじめに入力数字列の先頭から固定長のサイズ分数字を取り出し、次に入力数字列の2番目の数字を先頭とし、固定長のサイズ分数字を取り出す。以降、3番目、4番目…末尾まで各数字を先頭としながら固定長で数字を取り出していくことで入力数字列を分割する。例では入力数字列である41210213139を6数字の固定長で分割している。

分割が終わったらニューラルネットワークによる変換を行う。ニューラルネットワークの入力に対して、分割された数字列と、その数字列の前に出現した4数字を係り受け情報として与える。その入力を使いニューラルネットワークから「日本語の文字コード」(以下、日本語コード)または「文字ではないことを示すコード」(以下、非文字コード)のいずれか一方の出力を受け取る。例では日本語コードの「大会」、「を」、「開催」、「する」と、非文字コードの「FFFF」にそれぞれ変換している。

最後にニューラルネットワークの出力結果である日本

[†]岩手県立大学, Iwate Prefectural University

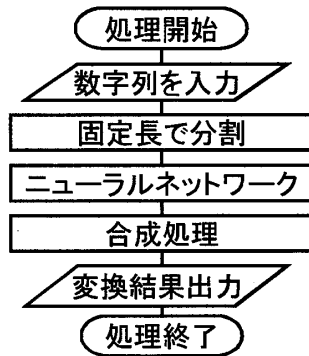


図2: 本手法の変換プロセス

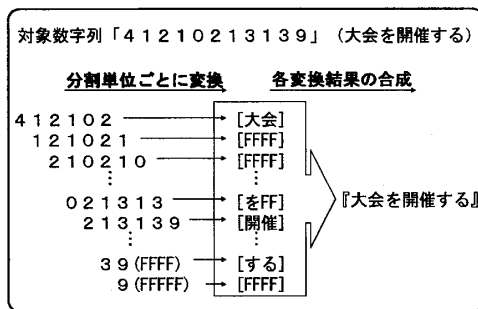


図3: 本手法における数字漢字変換

語文字コードと非文字コードを使い、合成を行うことで、入力された数字列の変換を完了とする。例では変換結果から文章を合成し「大会を開催する」を出力している。

このように、前処理として固定長で分割することで、どのような長さ、区切りで入力された数字列に対しても対応することが可能となる。また、固定長のサイズによって入力ノードの数を定めることが出来るため、ニューラルネットワークの規模を自由に定めることができる。

4. 評価実験

4.1 実験データおよび実験方法

入力数字列を分割し変換するまでの評価実験を、オープンデータを用いて行った。学習データには日本語の論文 [3] である、21,576 字分のテキストデータを用いた。このテキストデータを茶筌 [4] で形態素解析し、入力数字列と日本語コードを作成した。そして、作成した入力数字列を固定長サイズ 9 で分割したものに、分割した数字列の前方に出現した 4 数字と、正解データを付加し、1つの学習データとした。今回は、32,433 個の学習データが生成され、このうち 28,999 個を学習に使い、残りの 3,434 個で変換を行った。実験に使ったニューラルネットワークは、入力ノード数 52 個、中間ノード数 144 個、出力ノード数 144 個で構成しており、学習モデルには誤差逆伝播法を用いた。

4.2 実験結果および考察

実験結果を表 1 に示す。ノード単位での正解率は日本語コードが 96.4[%]、非文字コードが 99.3[%]、合わせて 98.2[%] であった。1つの分割数字列に対して完全な

表 1: 変換正解率

	日本語コード	非文字コード	総計
ノード単位	96.4[%]	99.3[%]	98.2[%]
分割単位	14.4[%]	81.4[%]	56.3[%]

表 2: 144 ビット中の平均誤りビット数

日本語コード	非文字コード	総計
5.1	1.0	2.6

形で変換に成功したことを表す、分割単位での正解率は日本語コードが 14.4[%]、非文字コードが 81.4[%]、合わせて 56.3[%] であった。先行研究のクローズドデータを用いた実験では、ノード単位で 97.3[%]、単語単位で 62.0[%] であった。クローズドデータよりも変換が難しいと考えられる、オープンデータを用いた実験においても、ノード単位では先行研究以上の変換精度を得ることができ、分割単位の変換精度も大きく劣ることはなかった。よって、この実験結果から本手法の有効性を示すことができた。

日本語コードの分割単位における変換率が非文字コードのそれより低くなっているが、これは表 2 から分かる通り、平均 5.1 ビットという非常に小さな誤りによるものである。今回の出力ノード数では、最大 9 文字までの日本語コード変換に対応しているが、実際には 9 文字の日本語コードは稀で、出力ノード全体を使う必要はない。よって、不要部分の情報を捨てることで、誤りビット数を減らすことが可能であると考えられる。また、既存の誤り訂正手法等を適用することで、さらに精度を向上させることが可能であると考えられる。

5. おわりに

本稿では、入力された数字列に対して、分割処理を施した後に学習を行うことで、一定規模のニューラルネットワークを利用しながらも、どのような長さ、区切りの数字列にでも数字漢字変換が可能である手法を提案し、オープンデータによる評価実験を行った。評価実験での変換精度から、本手法の有効性を示すことができた。

しかし、日本語コードの分割単位における変換精度は、平均して 5.1 ビットという小さな誤りによって低くなっている。よって、今後は、分割単位での変換精度を向上させるために、出力ノードの結果から不要部分の情報を除く処理や、誤り訂正機能などについて検討していく予定である。

参考文献

- [1] 鎌田竜也, 松原雅文, Goutam Chakraborty, 馬淵浩司: ニューラルネットワークを用いた携帯端末向け日本語入力手法における単語変換精度, 情報処理学会第 67 回全国大会講演論文集, 2J-4, pp.83-84 (2005)
- [2] 松原雅文, 荒木健治, 桃内佳雄, 枋内香次: 文字情報縮退方式を用いた帰納的学習によるべた書き文の数字漢字変換手法の有効性について, 電子情報通信学会論文誌 D-II, J83-D-II, No.2, pp.690-702, (2000)
- [3] 川嶋宏彰, 松山隆司: 連続状態モデル間の相互作用に基づく多視点動作認識, 情報通信学会論文誌 D-II, J85-D-II, No.12, pp.1801-1812, (2002)
- [4] 奈良先端科学技術大学院大学 情報科学研究科 自然言語処理学講座: ChaSen - 形態素解析器, <http://chasen-legacy.sourceforge.jp/>