

日本語に適した単語の誤入力訂正法とその大語い 単語音声認識への応用†

栗田 泰市郎†† 相沢 輝 昭††

単語音声認識や OCR を利用して文章や単語を計算機に入力する場合、文字ないし音節の置換誤りを生じる。本論文は日本語の単語入力において、単語辞書と2字組の頻度表などを用いて置換誤りを訂正する方法を提案する。さらにこの訂正法を単語音声認識と組み合わせて大語い単語音声認識に応用した例について述べる。本文では、まず誤り訂正における日本語と英語の違いについて検討し、日本語の誤り訂正に関する基本的な指針を得る。次に、従来の訂正法に誤字の2段階検出などいくつかの改善を加えた訂正法を提案する。基礎的な誤り訂正実験を行ったところ、本訂正法の有効性、とくに2段階検出の効果が確かめられた。模擬単語音声認識部と本訂正法による単語認識部(誤り訂正部)を用いての大語い単語音声認識シミュレーションでも良好な単語認識性能が得られた。たとえば、模擬単語音声認識部の単語音声認識率が95%のとき、単語認識率は約96%、単語認識部での処理時間は約0.3秒であった。本訂正法は他の訂正法の一つである、訂正に際して辞書の全単語を参照する方法に比べ、わずかに低い認識率を数十倍速い時間で実現できている。

1. ま え が き

単語音声認識による音声入力や OCR 入力は文章や単語を計算機に入力する有力な手段であるが、入力の際に誤りを生じるため、誤りの訂正などを行ってトータルの認識率を高める研究が古くからなされてきた。また、近年ワードプロセッサの普及に伴い日本語音声入力ワードプロセッサの要望が高まっている。その手段として、かな1文字が1音節にほぼ対応するという日本語の性質を利用して単語認識がよく用いられているが¹²⁾、ここでも認識精度を高めるため音節の認識誤りを訂正する方法が必要とされ、種々試みられている。

これらの文書あるいは単語データの入力における誤りは置換誤り(ある文字を他の文字にまちがえる誤り)には限られることが多く、単語内の置換誤りを訂正すれば単語としての認識率を高めることができる。そこで、単語辞書や文字の接続情報などの言語情報を利用した誤り訂正がおもに英単語に関して研究されてきた¹¹⁾⁻⁸⁾。

ここで見られた手法は二つに分類できる^{11), 4)}。一つは辞書(または単語の見出し表)の探索のくり返し、あるいは辞書に登録されている全単語の参照を基本とする方法である。他の一つは言語において文字の接続

のしかたが一様でなく、文字の組合せによって文字組みの出現頻度に大きな差があることを利用する方法である。前者を辞書法、後者を統計法と呼ぶことにする。統計法は言語における文字出現のマルコフ性を利用しているともいえるのでマルコフ法と呼ばれることもある^{6), 8)}。これらを利用した単語認識システムにおいては、誤り訂正の能力が最終的な単語認識能力に直接大きな影響を与えている。

一般に、辞書法は高い訂正率を得ているが、訂正の対象となる候補がかなり多く、訂正速度は遅い^{4), 6)}。統計法は、辞書を利用する場合もあるが³⁾、候補の数は少なく訂正速度は速い。しかし、一般に訂正率は低い。

ところで、近年の状況では音声認識や OCR などのデータ入力機器はオンライン的に用いられることが多く、認識の高速化はより重要な課題となっている。単語音声認識や OCR を誤り訂正と組み合わせて大語いの単語認識システムを構成できるが、処理の高速性という点では前記の統計法が有利である。しかしながら日本語の単語認識においてこの立場がとられたことは少ないようである¹¹⁾。OCR における阿部らの報告¹⁰⁾、音声認識における古井¹³⁾、外川ら¹⁴⁾、中本ら¹⁵⁾の報告がみられるが、これらの手法はいずれも辞書法に分類されよう。

本論文ではこのような観点から、統計法の立場にたって日本語に適した誤り訂正法を提案する。さらに、これを単語音声認識と組み合わせて大語い単語音声認識に応用する。本訂正法はほかにもかな文字を読み取る OCR などと組み合わせることができる。以下

† A Method for Correcting Errors on Japanese Words Input and Its Application to Spoken Word Recognition with Large Vocabulary by TAICHIRO KURITA and TERUAKI AIZAWA (Information Processing Research Division, NHK Technical Research Laboratories).

†† NHK 総合技術研究所情報処理研究部

2, 3 章では単語の構成要素をかな文字として論じ, 4 章では音節として論じる。

まず 2 章では誤り訂正における日本語と英語の違いをおのおのの単語の統計にもとづいて検討し, 日本語の誤り訂正に関する基本的指針を得る。

3 章では, 従来の訂正法に誤字の 2 段階検出などいくつかの改善を加えた誤り訂正法を提案する。かな鍵盤入力における打鍵誤りを題材に簡単な訂正実験を行い, 提案した訂正法の効果を確認する。

本訂正法は単語辞書, 2 字組 (digram) の頻度表, 入力時の置換誤りに関する confusion matrix を利用して訂正を行うが, これらのデータを記憶する容量は小さい方が好ましいのは当然である。3 章ではデータ容量の圧縮についても簡単にふれる。

4 章では実際の音声認識装置から得られた認識データをもとに模擬単音節認識部を作成し, それと本訂正法による誤り訂正部 (以下では単語認識部と呼ぶ) とにより大語い単語音声認識システムを構成する。このシステムの能力をシミュレーションによって求める。さらにこの結果を辞書法⁸⁾による結果と比較検討する。本訂正法による単語認識部は辞書法によるものと比べ, わずか低い単語認識率を数十倍速い速度で実現できることが明らかになる。

2. 誤り訂正における日本語と英語の違い

統計法による誤り訂正では, データとして n 字頻度表を利用する。 n 個の文字を組にしたものを n 字組 (n -gram) と呼び¹⁾, その種類は文字の種類数を m とすると m^n だけある。 n 字頻度表はこれらの n 字組が一般的な文章, あるいは辞書の見出しに現われる頻度を表にまとめたものである。

日本語の 6,460 語からなる辞書* から 2 文字以下の単語を除いた 5,704 語の辞書を編集し, そのかな見出

表 1 日本語と英語における n 字頻度表の例
Table 1 Examples of n -gram for Japanese and English words.

言 語	日 本 語		英 語 ¹⁾	
辞書の登録語数	5,704		2,755	
見出しの総文字数	24,558		16,530	
文字の種類 m	58		26	
n	2	3	2	3
記憶量*	3,364	195,112	676	17,576
ゼロ頻度率	59%	98%	約35%	約85%

* n 字頻度表が記憶すべき頻度の数 ($=m^n$)。

* ニュース文のかな漢字変換実験用の辞書¹⁰⁾。

しに関して n 字頻度表 (以下, 頻度表と略す) を作成した。このとき, ある文字組が見出しに現われる回数をその文字組の頻度とした。データを表 1 に示す。ここで, 濁点, 半濁点は 1 字として扱っており $m=58$ である*。このとき, 訂正能力を高めるため, 頻度表を単語長ごとあるいは n 字組が現われる文字位置ごとに作成するという方法⁷⁾ もあるが, 表 1 においては単語長にも文字位置にも依存しない方法をとった。比較した英語の例は Riseman らによるもので, 2,755 語の 6 文字単語が対象となっている⁴⁾。

表 1 において, 「ゼロ頻度率」は頻度表に記憶される m^n 個の頻度のうち, 値がゼロのもの占める割合である。この数値は誤り訂正において重要な意味をもつ。統計法による置換誤りの訂正は, 単語内の誤字を探し出す操作 (誤字検出) と, 誤字を正しい可能性の高い字に置き換える操作 (誤字訂正) から成る。これらの過程において, ゼロ頻度率が高いほど誤字を検出できる可能性が高く, 誤字訂正において正しい字の候補を絞ることができる。

表 1 では辞書の登録語数が異なるので単純な比較はできないが, 日本語は英語に比べて同じ n ならばゼロ頻度率が高いといえる**。 m が大きいためであろう。とくに $n=3$ の場合は 98% という高い値が得られている。しかし, 頻度表の記憶量も日本語のほうが大きく, $n=3$ の場合は英語に比べ約 11 倍となる。 $n=2$ の場合, 日本語の例はゼロ頻度率, 記憶量とも英語の $n=2$ と $n=3$ の中間の値となっている。

これらの結果から, 日本語において 3 字頻度表 (trigram) を利用した訂正は, 訂正能力は高いが, データの記憶容量, 訂正速度等の訂正にかかるコストもかなり高くなると予想される。Riseman⁴⁾ らは頻度表に記憶する頻度を 2 値化した binary n -gram を用いた。しかしこの結果, 約 99% という高い誤字検出力を維持しながら記憶容量が小さくなったが, 誤字訂正において数値的な処理ができなくなった。誤りの訂正率は, 2,755 語の辞書をベースにして最大で 62% と低い値になっている。

われわれは, 2 字頻度表 (digram) を利用し, 誤り訂正のアルゴリズムを強化することによってコストが小さくかつ訂正率の高い訂正法をめざした。もう一つの方向として, $n=3$ のゼロ頻度率 98% に着目し, 頻

* JIS かな符号の文字集合に合わせた。

** ゼロ頻度率は登録語数が多くなるにつれ緩やかに減少する⁴⁾。暗号学によれば, 登録語数よりも見出しの総文字数のほうがより直接的に影響するようである⁹⁾。

表 2 単語長の統計
Table 2 Statistics of word length for Japanese and English words.

	例 1 ^{*)}	例 2 ^{*)}	例 3 ^{*)}
言語	日本語*	英	語
語長	語数		
1	139	0	22
2	617	108	89
3	1,288	247	397
4	2,209	467	1,007
5	1,531	571	1,384
6	555	667	1,773
7	96	756	1,886
8	22	624	1,666
9	3	516	1,364
10以上	0	879	2,061
計	6,460	4,835	11,603
平均語長	4.0	7.0	7.3

* 語長はかな文字を単位として数えたもの。このとき濁点、半濁点は1字として扱っている。

度がゼロでない3字組のみを記憶するという方法も考えられる。この場合、3字頻度表へのアクセス時間を増加させないようなデータ構造にすることが重要である。

単音節認識やOCRの後処理として誤り訂正を使用する場合、単語の長さも考慮すべきことである。文字読取り部(OCR)と単語認識部(誤り訂正部)からなる単語認識システムを例にとって説明する。読取り部では各文字を一定の精度で読み取り、その結果を単語ごとに出力する。したがって読取り部分から出力された単語は何文字かの誤字を含んでいる可能性があり、それらを訂正する必要がある。このとき、単語長が短いほど1単語内で発生する誤字の数は少ない。すなわち、訂正すべき文字数が少ない。

表2は日本語と英語の辞書における単語長の統計である。少ない例ではあるが、英語より日本語のほうが平均語長が短いようである。もし読取り部の文字認識率が日本語でも英語でも90%であれば、平均的な4文字の日本語では1単語内で1字誤る確率は、 $0.9^3 \times 0.1 \times 4 \approx 0.29$ 、2字以上誤る確率は $1 - 0.29 - 0.9^4 \approx 0.05$ である。一方、平均的な7文字の英語では1字誤る確率は約0.37、2字以上誤る確率は約0.15である。この例で、日本語では1単語内1字の訂正を行うだけで約95%までの単語認識率が期待できるのに対して、英語は約85%しか期待できない。このように日

本語の誤り訂正では単語の平均語長が比較的短いため、1単語中1字の訂正能力をもつだけで単語認識システム全体としてかなりの能力が期待できる。

日本語は単音節認識により単語や文章を入力することが可能であるが、単語の平均的な音節数はやはり4程度であり、上と同様な能力が期待できる。このことは後に4章でもふれる。

なお、表1、表2で用いた日本語の辞書は後の実験すべてに使用するので、ここでその性質を調べておく。まず、この辞書のエントロピは3字組*i, j, k*の頻度分布を $p(i, j, k)$ などと表すと

単文字に対して

$$F_1 = -\sum_i p(i) \log_2 p(i) = 4.95$$

2字組に対して

$$F_2 = -\sum_{ij} p(i, j) \log_2 p(i, j) = 3.59$$

3字組に対して

$$F_3 = -\sum_{ijk} p(i, j, k) \log_2 p(i, j, k) = 2.34$$

(bits per letter)

であり、日本語4万字に対する今榮²¹⁾の結果と同様な値を示している。Shannon²²⁾による英文のエントロピに比べると F_1 は約0.8 bit 大きく F_2 は同程度である。また F_3 は約1 bit 小さい。このことは文字を予測するに際し、文字の種類は英文より2倍以上多いにもかかわらず2字組を用いれば予測の候補文字数を英文と同程度に絞ることが可能である、すなわち誤り訂正には有利であることを意味している。3字組の使用はさらに有利であろう。

次に、同じ語長の単語を比較したときそれらの間の異なり文字数 d を単語間距離と定義すると、語数 N の辞書において距離 d にある単語対の数 N_d の N に対する比 $n_d = N_d/N$ は、一つの単語に対しそれと距離 d にある単語の平均個数である¹⁰⁾。この n_d を前記辞書の3文字以上の単語に関して求めると、

$$n_1 = 18,978/5,704 \approx 3.3$$

$$n_2 = 119,685/5,704 \approx 35.0$$

$$n_3 = 1,336,239/5,704 \approx 234.3$$

となった。これらの数が大きいことは誤り訂正において後に述べる誤訂正が多くなりやすいことにつながり好ましくない。本辞書は語数 N が大きいため阿部¹⁰⁾や中本¹⁵⁾が実験で用いたものよりも n_d が大きく、きびしい条件になっている。

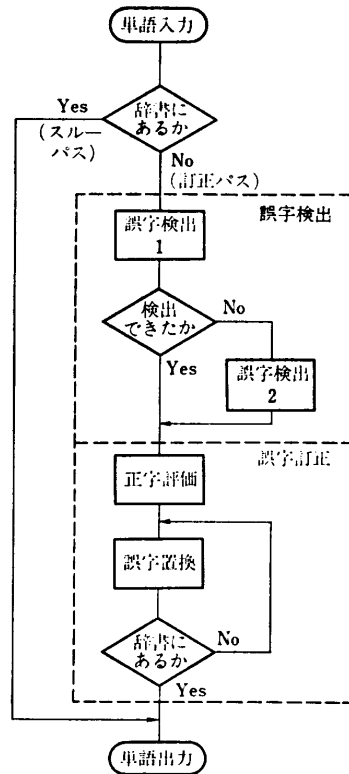


図 1 誤り訂正の流れ

Fig. 1 Flow of the correction of a substitution error per word.

3. 誤り訂正法とその能力^{17), 18)}

1 単語中 1 文字の置換誤りを訂正する方法について述べる。訂正されるべき単語は単語辞書に登録されているとする。訂正処理の大まかな流れを図 1 に示す。

単語が入力されるとまず辞書を探索する。単語が辞書に登録されていればそれは正しいと判断され、図 1 のスループスを通してそのまま出力される。登録されていないならば、訂正パスに入り誤り訂正がなされる。このように訂正に先だって辞書を探索する方法は Cornew³⁾ にも見られる。

3.1 訂正アルゴリズム

訂正パスは誤字検出と誤字訂正の二つの操作から成る。誤り訂正に使用するデータは 2 字頻度表および音節ないし文字を入力する際の confusion matrix である。confusion matrix はある文字をある文字に認識する m^2 個の頻度を記憶している。

(1) 誤字検出

誤りのある単語内から頻度がゼロである n 字組を見出すことは誤字を検出する一つの方法である⁴⁾。 n 字頻度表は辞書の見出しに関して作成されているの

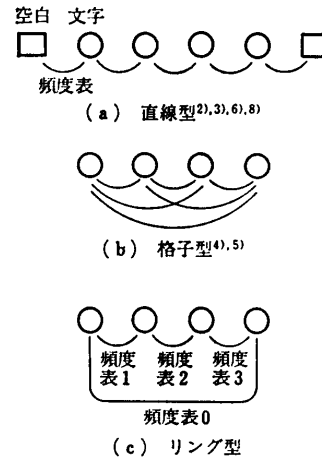
図 2 2 字頻度表の適用法 ($L=4$ の例)

Fig. 2 Application methods of digram (example of word length=4).

で、ある n 字組の出現頻度がゼロであれば、 n 字組は辞書の見出しに現われない、つまり、その n 字組のなかに誤字があることがわかる。これは一つの確実な誤字検出法である。この方法を以下、ゼロによる検出と呼ぶ。一つの頻度ゼロで n 個の誤字 (の候補) が検出される。

図 1 における「誤字検出 1」はゼロによる検出を行っている。しかし、3.2 節の実験結果で明らかになるように、2 字頻度表を利用したゼロによる検出では誤字を検出できる率が低い。そこで、誤字検出 1 の後に誤字が検出できたかどうかをチェックし、できていなければ「誤字検出 2」を行う。この構成を誤字の 2 段階検出と呼ぶことにする。

誤字検出 2 では、2 字組の頻度値にしきい値 D_{th} を設け、頻度が D_{th} 以下ならばその 2 字組のなかに誤字があると判断する。 $D_{th}=0$ ではゼロによる検出と等価になる。 D_{th} を変化させたときの状況は 3.2 節で論じられる。

ところで、2 字頻度表を単語に適用するに当たり図 2(a), (b) に示すような方法が試みられている。(a) は単語の前後に空白 (ブランク) を仮定し、頻度表を直線的に適用する方法であり、2 字頻度表をたとえば 1 文字目と 2 文字目の頻度表というように文字位置ごとに用意するとして、 $L+1$ 枚の頻度表が必要である (L は単語長)。(b) は頻度表をすべての文字の組合せに適用する方法で、 $L C_2$ 枚の頻度表が必要となる。図ではこれらを仮に直線型、格子型と呼ぶ。

(b) は強力であるが頻度表の記憶量が多い。たとえば、われわれの辞書には 9 文字の単語まで含まれて

いるので、 $C_2=36$ 枚の頻度表が必要となる。(a)は誤字の検出能力が低い。とくに、単語の語尾と語頭に誤字があるときに検出能力が低いようである。語尾と語頭の2字組の片側はつねに空白であり、情報量が不足するのではないだろうか。

ここでは図2(c)に示す適用法を用いた。単語の語尾と語頭に関する頻度表(図の頻度表0)を用意し、リング型に適用する。必要な頻度表の枚数は L 枚である。これにより、語頭や語尾も他の文字位置とまったく同様な処理が可能になる。これは次の誤字訂正においても有利である。確実な検出である誤字検出1で検出できる率は語尾において8%程度増加する。

(2) 誤字訂正

まず、ある文字が誤字に対してもとの正しい字(以下、正字と呼ぶ)である可能性を評価する。図1の「正字評価」である。誤字検出で検出された誤字が e 個あるとし、文字の種類を m とすると、1単語中に真の誤字は1字しかないと仮定しているため、評価すべき組合せは、 $e \times (m-1)$ 通りとなる。

評価値 g は次式で表されるものを用いた。頻度表の適用法はさきのリング型を用いている。

$$g_{ij} = K_{ij}^2 \times f_{1ij} \times f_{2ij} \quad (1)$$

ここで、文字の集合を $\{C_1, C_2, \dots, C_m\}$,

入力された単語を $W = X_1 X_2 \dots X_L$,

検出された誤字の位置を $p_i (1 \leq i \leq e)$ とすると、

検出された誤字は X_{p_i} で表せる。

X_{p_i} に対する正字の候補: $C_j (1 \leq j \leq m, \text{ただし}, X_{p_i} \neq C_j)$

g_{ij} : C_j が X_{p_i} の正字である可能性の評価

K_{ij} : 音節ないし文字を入力する際に C_j を X_{p_i} に誤る頻度

f_{1ij} : 2字組 $\begin{cases} X_{p_i-1} C_j (p_i \neq 1) \text{の頻度} \\ X_L C_j (p_i = 1) \end{cases}$

f_{2ij} : 2字組 $\begin{cases} C_j X_{p_i+1} (p_i \neq L) \text{の頻度} \\ C_j X_1 (p_i = L) \end{cases}$

g_{ij} をすべての組合せについて求め、 g_{ij} の大きい順を正字の候補(以下、正字候補と呼ぶ)の順位とする。 K_{ij}, f_{1ij}, f_{2ij} のいずれかがゼロであることにより、多くの文字が候補から除外される。 f_{1ij}, f_{2ij} に関してリング型の適用法を用いているため、語尾や語頭の文字でも他の文字位置と同程度に正字候補の数を絞ることができる。

図1の「誤字置換」では、第1の正字候補から順に誤字と置換し、そのたびに辞書を探索する。置換の結

果できた単語が辞書に登録されていればその単語を訂正結果として出力する。これを登録されている単語が生成できるまでくり返して誤りを訂正する。

(1)式において、 K_{ij} は confusion matrix から得ているが、 K_{ij} の2乗としたのは、(1)式における2字頻度表と confusion matrix の相対的な影響力を同等にするためである。われわれのある実験では K_{ij} を用いたとき誤りの訂正率86.7%に対し K_{ij}^2 を用いたとき訂正率88.7%となって K_{ij}^2 のほうが若干よかった。 K_{ij}^3 以上では K_{ij}^2 と大差ない結果を得ている。

正字候補の順位について、Cornew³⁾は誤字の位置に関してまず順位をつけ、その範囲内で正字候補の評価値によって順位をつけた。つまり誤字の位置と評価値で2段階に順位をつけている。これに対し本訂正法は評価値のみで順位をつける。両者を実際に比較したところ、本方法のほうが訂正率は数%高く、訂正速度は4割程度速かった。

以上の訂正アルゴリズムの形式的記述を付録1に示す。

3.2 訂正実験

JIS 配列のかな鍵盤から単語を入力する際の誤り、つまりキーの押しまちがいを対象に誤り訂正実験を行った。よくいわれるように、鍵盤入力においては置換誤りは誤りの一部にしかすぎない。しかし、本訂正法の訂正能力を知る基礎的な実験であるので、身近な鍵盤入力の置換誤りを対象とした。

単語辞書は表1のときに用いた5,704語のものを使用した。2字頻度表は辞書のかな見出しから図2(c)の方法に従って9枚作成した^{*}。このとき9枚の頻度表のゼロ頻度率はそれぞれ70%以上となっている。

confusion matrix は次のように作成した。通常、打鍵時には隣接キーにまちがえやすい。そこで単純に、あるキーの回りのキーに誤る頻度を1、他は0とした。

実験は、単語を辞書から選び出し、confusion matrix の作成と同じ基準で選んだ誤字1字を単語内の1字と置き換えて訂正プログラムに入力する。このような試行を5,704語の中から乱数によって選んだ単語100語に対し^{**}、可能な限り誤字と誤字位置を変えて

^{*} 辞書に登録されている最も長い単語が9字である。

^{**} このときの単語100語の集合におけるエントロピーは一例として

$F_1=4.9, F_2=2.4, F_3=3.1$ (bits per letter)

であった。また単語間距離に関する統計は一例として

$n_1=0.1, n_2=0.6, n_3=4.0$

であった。

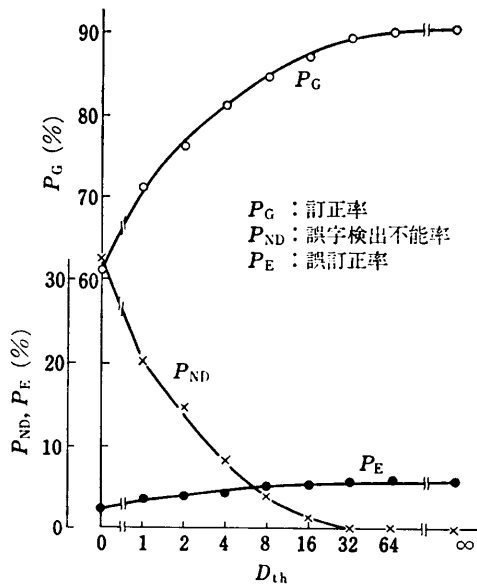


図3 2段階検出による訂正率の改善

Fig. 3 Improvement of correction rate by the 2-step error detection (increasing the threshold D_{th} of digram).

表3 誤り訂正の例

Table 3 Examples of spelling correction for Japanese words.

状 態	もとの単語	誤った単語	訂正結果
訂正可能	カイトウ	カテトウ	カイトウ
誤字検出不能	クミチガイ	クミチガテ	—
誤訂正	カイトウ	カイテウ	カイテイ
無訂正	テンラン	テンラク	テンラク

合計約 2,000 回行った。

結果を要約すると、単語認識部の出力は表3に例を示したような四つの状態に分類でき、その割合は D_{th} の増加に従って図3のように変化する。

表3において、誤字検出不能は誤字をまったく検出できなかったため単語がリジェクトされたことを示す。誤訂正はもとの単語とは異なる別の「正しい」(辞書に登録されている)単語を出力したことを示す。無訂正は誤字を含んでいる状態がすでに「正しい」単語となっているので図1のスループスを通して入力があるまま出力されたことを示す。訂正プログラムの訂正率は訂正可能が全体に占める割合である。

図3において、 $D_{th}=0$ の状態、つまりゼロによる検出のみで誤字検出を行った場合、訂正率は 61.3% とかなり低い値になる。訂正可能以外の大部分は 32.5% の誤字検出不能であり、 D_{th} の増加により訂正率の改善が期待できる。

表4 $D_{th}=\infty$ のときの訂正結果
Table 4 Correction rate and other results for $D_{th}=\infty$.

訂正可能	90.2%
誤字検出不能	0
誤訂正	6.1%
無訂正	3.7%
辞書探索数	2.6

誤訂正は最も好ましくない状態といえる。表3の例のように訂正結果から誤りがわかりにくくなることが多いためである。 D_{th} を増加させる際に心配されたことは、誤字検出不能が減るかわりに誤訂正が増すことであった。しかし、図3に見るように誤字検出不能の急速な減少は大部分が訂正可能にかわり誤訂正の増加は少なかった。

$D_{th}=\infty$ のとき訂正率は最高となった。表4はこのときの各状態の割合である。訂正率は 90.2% と高い値になっている。 $D_{th}=\infty$ ということは、ゼロによる検出ができなかった場合は単語内のすべての文字に誤字の可能性があると判断する、ということである。

表4の辞書探索数は辞書を探索する回数の1試行当りの平均であり、訂正速度の一つの目安である。図1の最初の探索があるため、この値は必ず1以上になる。辞書探索数が小さいことは訂正速度がはやいことにつながり、それはまた(1)式の g_{ij} の有効性の一つの目安でもある。本実験といろいろ条件は異なるが、参考として、Cornew³⁾は約70%の訂正率を約10回の辞書探索で得ている。本実験では2.6回という少ない辞書探索で約90%の訂正が可能となっている。

なお、表4の結果は辞書から単語を選ぶ乱数を変えてもほとんど変わらない。たとえば、ある別の乱数のとき訂正率 90.6%、誤訂正 5.9%、無訂正 3.6%であった。

3.3 2字頻度表の圧縮

前節では文字位置ごとに9枚の2字頻度表を使用しており、各頻度表のゼロ頻度率は70%以上であった。図2(c)にならってこれらを頻度表0~8と呼ぶ。一つの頻度は2バイトで記憶しており、頻度表全部の記憶容量は $58^2 \times 9 \times 2 = 59.1$ kバイトである。

この量はやや大きいので頻度表0、および頻度表1~8を1枚に圧縮した頻度表1'の2枚を使用することにし(つまり、頻度表0と1'の2枚で図2(c)のリングを構成する)、また一つの頻度を1バイトに圧

縮して*, 訂正実験を行った¹⁸⁾. この状態で頻度表の記憶容量は $58^2 \times 2 \times 1 \approx 6.6$ k バイトとかなり小さくなっている. このように圧縮しても訂正能力はほとんど変わらないことが報告されている¹⁸⁾.

ここで頻度表 1' は表 1 に日本語の $n=2$ の例として示したものと同一のものである. この頻度表のゼロ頻度率は表 1 から 59% であるが, このようにゼロ頻度率が低下しても訂正アルゴリズムの改善により訂正率の悪化は小さくなっている.

頻度表以外のデータについては, confusion matrix の記憶容量は 3.3k バイトであり頻度表と同等以下である. 辞書はその大きさと用途の広さから大容量の 2 次記憶に置くものと考えている.

文献 18) では, 上の圧縮と同時に誤字の発生条件を 3.2 節より複雑な条件に変えているが, これによっても訂正能力はあまり劣化していない.

4. 大語い単語音声認識への応用

4.1 模擬単音節認識部

この章では前章で述べた誤り訂正法を大語い単語音声認識に応用することを試みる. 以下では単音節認識部と, 音節認識の誤りを単語ごとに訂正して単語を認識する単語認識部からなる大語い単語認識システムを考える. 単音節認識部は音声信号を受けて音節の第 1 候補の列を単語ごとに出力し, 単語認識部はその列の中に認識誤りがあれば訂正して出力する. 図 1 の誤り訂正プログラムを単語認識部としてそのまま用いるので, 訂正できるのは 1 単語中 1 音のみである.

単音節認識部として実際の装置を用いるのではなく, 以下に述べるような模擬単音節認識部を用いる.

専用の単音節認識装置が入手できなかったので, 市販の単語音声認識装置に音節を一つの単語として認識させてみた. この装置は特定話者用のものである. 五十音と濁音, 半濁音を含む 68 音節を男性 10 人, 女性 10 人に 5 回ずつ発声してもらい, 装置の認識特性を調べた. この装置の単音節認識率は 62.2% であった.

このとき得られた confusion matrix の対角成分 (音節が正しく認識できることを示す成分) に定数をかけると, 音節の認識誤りの性質をそのままにして単音節認識率を任意に設定できる. これを模擬単音節認識部とした. 前記装置の子音に関する confusion matrix

* 頻度表 0 と頻度表 1' は各要素が 2 バイト整数である 58×58 の行列で表しうが, 値が 255 以上となる要素はきわめて少ない (頻度表 0 では 0, 1' では 6 個) ので, 255 以上の要素を 255 に置きかえ, 各要素を 1 バイト整数に圧縮した.

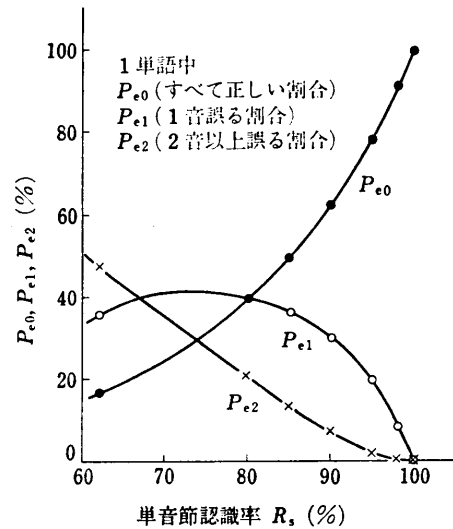


図 4 模擬単音節認識部の特性
Fig. 4 Performance of the simulated syllable recognition block vs. Japanese syllable recognition rate R_s .

を付録 2 に示す.

図 4 は模擬単音節認識部のみで単語認識シミュレーションを行った結果である. 以下では単音節認識率を R_s で表す. R_s は 68 音節を均等に発声した場合に音節認識結果の第一候補が正解である確率である. 図において, 実用的な $R_s=90\%$ 以上の領域では, 1 単語中 1 音誤る割合が誤りの大部分を占めている. $R_s=90\%$ のとき 1 音誤り (図の P_{e1}) は約 30%, 誤りがない場合 (P_{e0}) は約 63% であり, 1 音誤りを訂正することにより 93% までの単語認識率が期待できる. これは訂正する前の 63% に比べかなり大きな値である. $R_s=95\%$ では単語認識率 98% まで期待できる.

4.2 単語認識シミュレーション

上記の模擬単音節認識部と誤り訂正とにより大語い単語音声認識シミュレーションを行った.

3 章と同様に 6,460 語の辞書から 2 音節以下の単語を除いた 5,496 語の辞書を編集した. そして, この辞書の見出しの音節表現に関して 2 字頻度表を作成した. つまり, 68 音節あるので 2 音組の種類は 68^2 である. また, 頻度表は 3.3 節で述べたような圧縮された形で使用し, 記憶容量は $68^2 \times 2 \times 1 \approx 9.0$ k バイトとなっている. 誤り訂正に使用する confusion matrix は模擬単音節認識部のもを用いる.

5,496 語の辞書からランダムに, 1,000 語を選び模擬単音節認識部に入力する. 模擬単音節認識部では confusion matrix にもとづいて乱数により音節の第

表 5 単語認識シミュレーションの結果
Table 5 Result of the simulation of Japanese word recognition.

単音節認識率 R_s	90%	95%
P_{e0} (図 4 参照)	63%	78%
単語認識率	89%	96%
辞書探索数	1.7	1.5
全試行数	1,000	1,000
1音誤り	299	199
2音以上の誤り	73	19
訂正可能	260	175
訂正率*	87%	88%

* 1音誤りに対する訂正可能の割合

一候補の列を出力する。単語認識部はこれを受け、誤り訂正を行って単語を出力する。誤字の2段階検出における D_{th} は 3.2 節の結果から $D_{th} = \infty$ とした。

表 5 は $R_s = 90\%$ と $R_s = 95\%$ の場合の結果である。表の P_{e0} は図 4 に示したもので、単語認識部の出力における最終的な単語認識率はその下に示されている。

$R_s = 95\%$ では、1単語の認識当り 1.5 回と少ない辞書探索数にもかかわらず、単語認識率は 96% と高い値になっている。1音の置換誤りを訂正できた状態である「訂正可能」は 1,000 語中 175 語であり、1音誤り 199 語に対する訂正率は 88% であった。

4.3 辞書法との比較

提案した統計法による単語認識部の性能を辞書法によるものと比較する。

Shinghal は高い認識率が得られる DA (Dictionary Algorithm) 法のリストを示した⁸⁾。これは特別なアルゴリズムによる文字認識部と組み合わせて用いられているが、第一候補のみを出力する単音節認識部に適合させることは簡単である。この DA 法を図 1 の訂正パスに置いたプログラムを作成した。付録 3 にその形式的記述を示す。DA 法は誤り訂正に際し、辞書の全単語を参照している。

4.2 節と同様な方法でシミュレーションを行った。試行数は 500 単語とし、1単語当りの認識時間を CPU 時間により測定した。ここでいう認識時間は単語認識部における処理時間である。計算機は TI の DS 990 ミニコンである。提案した方法および DA 法による結果を表 6 に示す。表 6 における「認識率」は表 5 と同様に訂正後の最終的な単語認識率である。

表 6 から明らかなように、提案した方法は DA 法

表 6 提案した方法と DA 法⁸⁾ の比較
Table 6 Comparison of the proposed method vs. DA method.⁸⁾

R_s (%)	提案した方法		DA 法	
	認識率 (%)	認識時間 (sec)	認識率 (%)	認識時間 (sec)
80	75	0.48	83	31
85	81	0.39	86	26
90	88	0.36	92	19
95	97	0.29	98	10

表 7 $R_s = 90\%$ のときの結果の分析
Table 7 Analysis of the result of Table 6 for $R_s = 90\%$.

状 態	提案した方法		DA 法	
	1音誤り	2音以上の誤り	1音誤り	2音以上の誤り
訂正可能	128	0	125	22
誤訂正	15	8	18	15
無訂正	8	0	8	0
2音誤りを検出	0	5	0	0
訂正不能	0	24	0	0
訂正率 (%)	85	0	83	59

よりやや低い認識率を 30~60 倍速い認識速度で実現している。とくに $R_s = 95\%$ のときは両者の認識率の差は 1% 程度にすぎない。このとき提案した方法は 1単語を約 0.3 秒で処理できる。 $R_s = 90\%$ のときは約 4% の差である。また、単語を正しく認識できない場合のなかで最も好ましくない誤訂正も、提案した方法は DA 法と比べ同数以下である。この様子を表 7 に示す。

表 7 は 1音以上誤りがある場合に対する認識結果の分析である。表に示したのは $R_s = 90\%$ の場合であるが、他の場合でも同様な傾向を示す。この場合、500 単語中で模擬単音節認識部が 1単語につき 1~3音誤る割合はそれぞれ、151, 32, 5 語だった。4音以上誤ることはなかった。表の「2音誤りを検出」と「訂正不能」は、それぞれの理由により単語がリジェクトされた状態である。

表 7 をみると、提案した方法と DA 法の認識率の差は 2音以上誤りがある場合 (大部分は 2音誤りの場合) の訂正能力のみである。したがって、1音誤りが誤りの大部分を占めているとはいえ、図 1 の単語認識部に 2音誤りの訂正機能を加えれば、認識能力はさらに改善できよう。1音誤りに対する訂正率はわずかではあるが提案した方法のほうがまざっている。

DA法が2音誤りに対して訂正能力をもっているとはいえ、その訂正率は約60%と低い値である。この程度の訂正率は図1の誤り訂正アルゴリズムをそのまま2音誤りに拡張しても達成できるのではないかと考えられる。具体的には、誤字訂正において正字候補の評価値を誤字を1単語内に2字仮定してさまざまな組合せについて求め、評価する方法などが考えられる。この点は今後の検討課題である。

なお、表6の認識時間は R_c が小さくなると増加するが、これは R_c が小さくなると図1の訂正パスに入る割合が増えるためである。訂正パスに入った場合の処理時間の平均は R_c によらずほぼ一定で、提案した方法で約1秒、DA法で約50秒であった。

ところで、最近、文字の接続情報を認識された音節や文字の下位候補にまで適用し、日本語の単語認識率を向上させたという報告がなされている^{19),20)}。本方法においても、下位候補や音節間の距離データを利用すれば図1の誤字訂正において正字候補をより絞れることが期待でき、表7の誤訂正はさらに減るであろう。下位候補や距離データの利用も提案した方法の訂正能力ひいては単語認識能力を改善できると考えられる。

5. む す び

日本語を入力する際に生じる置換誤りの訂正法とその大語い単語音声認識への応用について述べた。

本訂正法は文字の接続情報を利用して訂正を行う統計法に分類されるものである。従来の方法にいくつかの改善を加えた結果、基礎的な訂正実験で高い訂正率が得られた。とくに、誤字検出における2段階検出の採用は、日本語の性質とあいまって、高い訂正率を記憶容量の小さい2字頻度表(digram)を用いて実現することを可能にした。

模擬単音節認識部を用いての大語い単語音声認識シミュレーションにおいては、単音節認識率が95%のときに、96%程度の単語認識率が認識時間約0.3秒で得られた。これは、辞書法による認識シミュレーションの結果に比べ、認識率は約1%低く認識速度は約30倍速い結果となっている。両者に認識率の差が出る原因を分析したところ、それは1単語内に誤りが2音ある場合の訂正能力であることがわかった。1単語内に誤りが1音のときの訂正率はほぼ等しい。

ここで報告した単語認識部(誤り訂正部)は単音節認識部からの情報として単音節の第一候補のみしか用

いていない。下位の候補とそれらの音節間距離データを利用すれば、より高い単語認識能力が期待できる。現在、これらの情報を利用して1単語内に誤りが2音ある場合まで訂正できる訂正法を検討している。

単語辞書を利用した単語認識システムの性能は辞書の登録語数が大きくなると低下することが知られている^{4),7),13),14)}。認識シミュレーションで用いた辞書は5,496語から成る。この大きさは過去の報告で用いられた辞書の大きさに比べて小さくはないといえよう。数万語から成る汎用辞書を用いて認識能力を測る必要がある一方で、その大きな辞書を分野別に分類するとすれば各分野別の辞書の大きさは数千語になるだろうという考えもある。辞書を分野別に分類することは、かな漢字変換など他の言語処理においても有用である。現在、NHK編新用字用語辞典をもとに分野別辞書の作成を進めている。

謝辞 おわりに、日ごろご指導いただくNHK総合技術研究所情報処理研究部の町田部長、沓沢主任研究員に深く感謝する。

参 考 文 献

- 1) 川合 慧: 英文綴り検査法, 情報処理, Vol. 24, No. 4, pp. 507-513 (1983).
- 2) Harmon, L. D.: Method and Apparatus for Correcting Errors in Mutilated Text, U.S. Patent, No. 3188609 (June 8, 1965).
- 3) Cornew, R. D.: A Statistical Method of Spelling Correction, *Inf. Control*, Vol. 12, pp. 79-93 (1968).
- 4) Riseman, E. M. and Hanson, A. R.: A Contextual Postprocessing System for Error Correction Using Binary n -Grams, *IEEE Trans. Comput.*, Vol. C-23, No. 5, pp. 480-493 (1974).
- 5) Hanson, A. R., Riseman, E. M. and Fisher, E.: Context in Word Recognition, *Pattern Recogn.*, Vol. 8, No. 1, pp. 35-45 (1976).
- 6) Shinghal, R. and Toussaint, G. T.: A Bottom-up and Top-down Approach to Using Context in Text Recognition, *Int. J. Man-Mach. Stud.*, Vol. 11, No. 2, pp. 201-212 (1979).
- 7) Hull, J. J. and Srihari, S. N.: Experiments in Text Recognition with Binary n -Gram and Viterbi Algorithms, *IEEE Trans. PAMI*, Vol. PAMI-4, No. 5, pp. 520-530 (1982).
- 8) Shinghal, R.: A Hybrid Algorithm for Contextual Text Recognition, *Pattern Recogn.*, Vol. 16, No. 2, pp. 261-267 (1983).
- 9) 加藤正隆: 基礎暗号学, 数理科学, No. 179, pp. 76-83 (1978).


```

度を表す)
{ $g_{ij}$  を大きい順にソートする};
{ $g_{ij}$  がいちばん大きい  $C_j$  と  $X_{pi}$  の組について  $X_{pi}$ 
を  $C_j$  に置換する};
while (W が辞書にない) do begin
  {前に置換した  $C_j$  を  $X_{pi}$  にもどす};
  { $g_{ij}$  が次に大きい  $C_j$  と  $X_{pi}$  の組について  $X_{pi}$  を  $C_j$ 
に置換する}
end;
Step 3: {単語 W を出力する};
Step 4: end. [訂正アルゴリズム終了]

```

付録 2

4.1 節で述べた市販の単語音声認識装置を単音節認識に用いた場合の子音に関する confusion matrix を付表 1 に示す。子音認識率は 62.6% であり、単音節認識率 62.2% に近い。すなわち音節認識誤りの原因はほとんど子音の誤りである。母音認識率は 98.5% だった。

付録 3

4.3 節における DA 法⁸⁾ の適用法を以下に示す。記法等は付録 1 と同様である。

Dictionary Algorithm

```

{単語  $W=X_1X_2\cdots X_L$  を入力する};
if (W が辞書にある) then go to Step 1;
for  $i=1$  to  $N$  do [ $N$  は辞書の単語数]
  if ( $W_i'$  の語長) $=L$  then begin
    [ $W_i'=X_{i1}'X_{i2}'\cdots X_{iL}'$  は辞書の  $i$  番目の単語]
     $g_i \leftarrow 1$ ;
    for  $j \leftarrow 1$  to  $L$  do
       $g_i \leftarrow K_{ij} \times f_{ij} \times g_i$ 
      [ $K_{ij}$  は  $X_{ij}'$  を  $X_j$  に認識する頻度,  $f_{ij}$  は 2 字
組  $X_{ij}'X_{i,j+1}'$  の頻度, ただし  $j=L$  のときは
 $X_{ij}'X_{i1}'$  の頻度]
    end;
  W ← (最も  $g_i$  の大きい  $W_i'$ );
Step 1: {単語 W を出力する};
end. [DA 法終了]

```

(昭和 58 年 11 月 11 日受付)

(昭和 59 年 4 月 17 日採録)