

漢字部首情報からの日本語単語の推定

梅田三千雄

大阪電気通信大学

漢字の部首などの部分構造情報をもとに、日本語単語辞書を用いることによって、どのくらいの単語が決定でき、漢字が推定できるかについて検討した。単語としての文字の性質を定量化するために、ここでは単語確定率と確定文字率の二つの評価尺度を定義した。まず、112種類の部首共通文字集合について、部首情報に基づく単語確定率を求め、ランダムに候補文字種を選択した場合と比較した。ついで、「偏」と「つくり」で構成される文字集合の性質を単語確定率と確定文字率により検討した。この結果、漢字の「偏」がわかると約72%の単語を決定でき、「つくり」がわかると約95%の単語を決定できることが明らかになった。

Estimation of Japanese Words
from the Knowledge about Kanji Radicals

Michio UMEIDA

Osaka Electro-Communication University
18-8 Hatsu-cho, Neyagawa-shi, Osaka 572, Japan

This paper discusses how many words can be determined and how many characters can be estimated by utilizing a Japanese word dictionary supposing only radical structures of kanji characters are correctly recognized. Two evaluation measures such as a word determination rate and a determinable character rate are defined to quantitatively evaluate some characteristics of each character as a word component. First, a word determination rate is calculated for character sets of common radicals and that is nearly equal to the rate for candidate characters selected at random. Next, two rates are calculated for the set of characters composed of left- and right-hand radicals. It is shown that about 72% of Kanji characters contained in Japanese words can be determined only from the structure about left-hand radicals and about 95% can be determined from that about right-hand radicals.

1. まえがき

我々が日常の情報伝達媒体として使用している漢字は、表意文字であるが故に文字種が非常に多く、その形も一般に複雑で、しかも形状の極めて類似した文字が多数存在する。また、日本語ワープロがかなり普及した現在でも、漢字は手書きされることが多く、手書きによる文字変形が大きいのも特徴である。これらの理由により、機械による漢字の認識、特に手書き漢字の高精度な認識は、これまでの活発な研究にもかかわらず、パターン認識における難しい課題の一つになっている[1][2]。

これまでに提案されている手書き漢字認識のアプローチには、大別すると、パターン整合法とストローク構造解析法の二つがある。パターン整合法とは、文字パターンそのもの、または文字パターンから抽出した特徴パターンを用いて、標準パターンとの整合処理により、文字を認識するものである。これに対して、ストローク構造解析法とは、文字パターンを構成するストロークセグメント群を抽出し、それらの種類や組み合わせ、配置関係などを手がかりに、文字を認識しようとするものである。

手書き漢字認識研究の初期のころには、ストローク構造解析法に基づく文字認識の試みがいくつか報告されている[3][4][5]。つまり、漢字のもつ構造に着目して、抽出したストローク群をもとに、部首に相当する部分構造を理解しながら文字認識を行おうとするものである。しかしながら、これらの研究では、漢字特有の文字種の多さや形状の複雑さの問題を克服することができず、漢字の構造記述や部分パターン抽出の困難さを指摘するに留まっている。

一方、パターン整合法では、文字種の多さは単に辞書規模や処理速度、認識精度などだけの問題であり、認識アルゴリズムの構築にはそれほど支障にはならない。まして、文字形状の複雑さは、むしろ認識のための特徴として積極的に利用することができる。このため、現在ではパターン整合法に基づく手書き漢字認識が主流となり、多くの研究により、かなり精度の高い認識手法も開発されている。しかし、パターン整合法では、文字全体を一つのパターンとして捉え、全体の形状として最もふさわしいと思われる文字種を認識結果とするため、ある文字の人力に対して、部首などの部分パターンの異なる、思いもつかぬ文字種が認識結果として出力されることがある。

これに対して、我々が漢字を認識し、理解する場合には、漢字のもつ部首などの部分構造の情報が大きな役割を果たしている。例えば、何か読めない漢字に出くわした時には、部首を手がかりに辞書を検索して、該当する漢字を見つけ出し、その読みや意味を調べる。また、ある漢字を誰かに伝えようとする時には、例えば「ごんべんににん、にんはにんたいのにんです。」

とか、あるいは「ごんべんににん、にんはやいばにこころです。」と言うことによって、「認」という漢字を相手に伝達することができる。

さらに、何かある読みづらい漢字に出くわした時には、その漢字の一部の情報をもとに、前後の文字を頼りにして、漢字を推測したり、単語としての意味を理解したりすることができる。また、漢字の一部を度忘れした時など、もっともらしく適当にごまかして書いても、単語あるいは文節、文章として、相手に正しく理解させることができることもある。人間は、このような柔軟な処理をすることによって、現在の文字認識アルゴリズムと比べて、はるかに高度で正確な文字の認識を行っている。

本報告では、文字認識における人間の高度で柔軟な処理のうち、日本語としての単語知識をもとに、漢字の部首情報のみを手がかりとした、単語の推定と漢字の決定について分析する。すなわち、日本語単語辞書を用いることにより、部首などの部分構造の情報しか明らかでない漢字に対して、それを含むどのくらいの単語が推定でき、その文字を決定できるか、つまり単語を構成する文字の単語としての性質を定量化することを試みる。

2. 評価尺度の設定

2.1 文字種と単語辞書

単語を構成する文字の単語としての性質を分析するためには、まず対象として扱う文字種を決定し、単語辞書を作成する必要がある。

対象とする文字種は、文字認識アルゴリズムの研究で頻繁に取り上げられる認識対象を参考に、常用漢字1945文字とひらがな75文字（清音字46、濁音字20、半濁音字5、小文字4「っ、ゃ、ゅ、ょ」）から成る合計2020文字とする。

一方、単語辞書は、学研国語大辞典[6]をもとに、親見出し語に対する正書法として記されているものの中から、上述の常用漢字とひらがなだけで構成されているものを選択して作成した。ただし、動詞や形容詞などの活用する語は語幹のみとし、ひらがなだけの単語や連語、慣用句、ことわざは除外した。さらに、単語長は2文字以上かつ6文字以下のものに限った。

こうして作成した日本語単語辞書の規模を表1に示す。単語総数は約52000語で、その約7割が2文字単語であり、平均の単語長は2.39文字となる。ただし、上述の文字種のうち、その文字を含む単語が全く存在しないものは「勺、朕、岬、匂、ぢ、ぶ」の6文字であった。

2.2 単語確定率

単語を構成する文字の単語としての性質を議論するため、まず単語確定率[7]なる評価尺度を設定する。

単語確定率 (Word determination rate) とは、対象とした単語辞書に含まれる全ての単語を構成する各文字のうち、その文字が正しく認識できなくても、単語を成す前後の文字がわかれば、単語としてその文字と混同する他の文字がなく、単語を一意に確定することができるものの割合で定義する。

この定義によれば、前述の日本語単語辞書における単語確定率は 27.45% になる。つまり、単語を構成する全ての文字のうち、その約 1/4 は前後の文字から単語を推定できるという性質をもっていることになる。しかし、これを 2 文字単語に限って見ると、その単語確定率は 0.31% と極めて小さく、ほとんどの単語を決定することができない。

つぎに、ある共通の性質をもつ文字だけを抽出した文字集合 W において、その集合を成す文字のみに着目した単語確定率を文字集合 W の単語確定率とする。さらに、ある文字「X」に着目して、それを含む単語のうち、その文字がわからなくても、前後の文字から単語を一意に決定できるものの割合を文字「X」の単語確定率と定義する。よって、単語確定率が 100% の文字は、単語を構成する前後の文字がわかれば、単語としての検定を行うことにより、それを含む全ての単語が一意に決まる文字である。また、逆に単語確定率が 0% の文字は、このままでは単語辞書を用いる効果が全くない文字である。しかし、何らかの方法によって、その文字に対する候補文字種が限定できれば、単語確定率は大きくなる可能性がある。

そこでいま、各文字に対する候補文字種を全字種のなかからランダムに選択して、対象を限定した場合の単語確定率を求めてみる。すなわち、文字種の総数を N 個とし、このなかから n 個の文字がラン

ダムに選択されるものとする。なお、文字「X」に対する候補文字種として、1 個は必ず文字「X」が選択され、残りの (n-1) 個の文字が (N-1) 個の文字種のなかからランダムに選択されるものとする。

ある単語中の文字「X」について、文字「X」以外にも単語を成す文字の個数を U とすると、単語を成さない文字の個数は N-U-1 (=V) となる。したがって、文字「X」以外の文字種が等確率でランダムに (n-1) 個選択されると、そのなかに単語を成す文字が p 個、単語を成さない文字が q 個含まれている確率 P(n) は、

$$P(n) = \frac{{}^u C_p \cdot {}^v C_q}{{}^{n-1} C_{p+q}} \quad (1)$$

ただし、U+V=N-1、p+q=n-1 で表される。ここで、

$$p=0, q=n-1 \quad (2)$$

の時には、単語中の文字「X」は一意に決定されることになる。したがって、その確率 Q(n) は、

$$Q(n) = \frac{{}^u C_{n-1}}{{}^{n-1} C_{n-1}} \quad (3)$$

で与えられる。この確率を全字種について求め、各文字種を含む単語数で重み付けすることにより、ランダムに候補文字種を選択した時の、選択文字数 n に対する単語確定率を導くことができる。

表 1 単語辞書の規模

単語長	単語数
2	36398
3	11570
4	3396
5	521
6	77
合計	51962

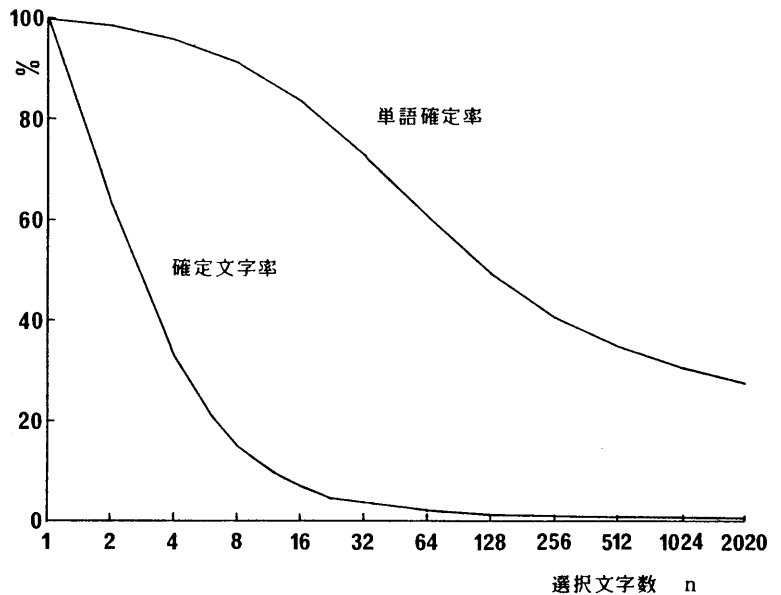


図 1 選択文字数に対する単語確定率と確定文字率

前述の単語辞書について、式(3)から算出した、選択文字数nに対する単語確定率を図1に示す。これより、何らかの方法によって、各文字に対する候補文字種をランダムに16個に絞ることができれば単語確定率は約83%に、さらに8個に絞ることができれば約91%にまで向上することがわかる。

2.3 確定文字率

単語を構成する文字の性質を議論する、もう一つの評価尺度として、確定文字率[7]を定義する。

確定文字率(Determinable character rate)とは単語確定率が100%になる文字種の全字種に対する割合、つまりある文字がわからなくても、単語を成す

他の文字から、その文字を含む全ての単語を一意に決定することができるという性質をもつ文字の割合とする。さらに、文字集合Wにおいて、その集合を成す文字のうち、単語確定率が100%になる文字の割合を文字集合Wの確定文字率と定義する。

前述の単語辞書において、単語確定率が100%になる文字はわずかに「峠、霧、ぜ、ば、び、ふ、へ、ほ、ゆ、よ、を」の11文字しかなく、その確定文字率は0.54%と極めて小さい。確定文字率も、何らかの方法によって、候補文字種を限定することができれば、大きくなる可能性がある。

ここでも、全字種からランダムにn個の候補文字種を選択して、対象を限定した時の確定文字率を求めて

表2 抽出した部首の種類とそれを含む漢字の個数

[偏]				
・汁(103)	・化(86)	・打(78)	言(56)	木(53)
系(52)	月(30)	金(28)	・防(28)	・忙(27)
土(26)	・行(25)	女(24)	・私(22)	口(19)
石(13)	貝(12)	日(12)	・犯(12)	・礼(12)
酉(11)	日(10)	車(10)	米(10)	・飢(8)
・初(8)	・距(8)	王(7)	馬(7)	弓(6)
山(6)	舟(6)	方(6)	目(5)	牛(5)
・次(5)	・列(5)			
[つくり]				
・刈(27)	・攻(25)	頁(19)	・役(13)	力(12)
欠(11)	・准(11)	・邦(11)	寸(11)	・印(9)
斤(7)	又(7)	・杉(7)	月(7)	寺(6)
反(6)	方(6)	・行(6)	各(5)	且(5)
見(5)	支(5)	尺(5)	主(5)	召(5)
青(5)	束(5)	皮(5)	包(5)	羊(5)
・悦(5)	・渴(5)	・限(5)	・孔(5)	・浅(5)
・場(5)	犬(5)	・化(5)		
[冠]				
・字(48)	・芋(42)	・笑(21)	・介(16)	田(14)
・文(13)	日(13)	雨(12)	・争(9)	・究(8)
・劳(6)	・堂(5)			
[脚]				
心(34)	口(23)	貝(22)	木(18)	・点(16)
土(16)	日(15)	・兄(14)	十(13)	女(9)
寸(9)	月(9)	田(9)	・六(8)	子(5)
[垂れ]				
・広(26)	・尼(13)	・病(12)	・厄(9)	・慮(6)
・戻(6)				
[構え]				
・囚(14)	門(12)			
[にょう]				
・込(50)	・赴(5)			

(・印は、その部首がJIS第1水準漢字にないため、その部首を含む最も簡単な漢字で表したものである)

みる。いま、ある文字「X」に着目し、その文字を含む全ての単語について、単語として文字「X」と混同する「X」以外の異なる文字種の総数をm個とする。選択されるn個の文字のうち、1個は必ず文字「X」が選択され、残りの(n-1)個の文字が(N-1)個の文字種のなかからランダムに選択されるものとする。すると、選択される文字のなかに文字「X」と混同するものがなく、全ての単語が一意に決定できる確率R(n)は、単語確定率の場合と同様に、

$$R(n) = \frac{N-m-1 C_{n-1}}{N-1 C_{n-1}} \quad (4)$$

で表される。したがって、これを対象とした全ての文字種に対して求めることにより、確定文字率を算出することができる。

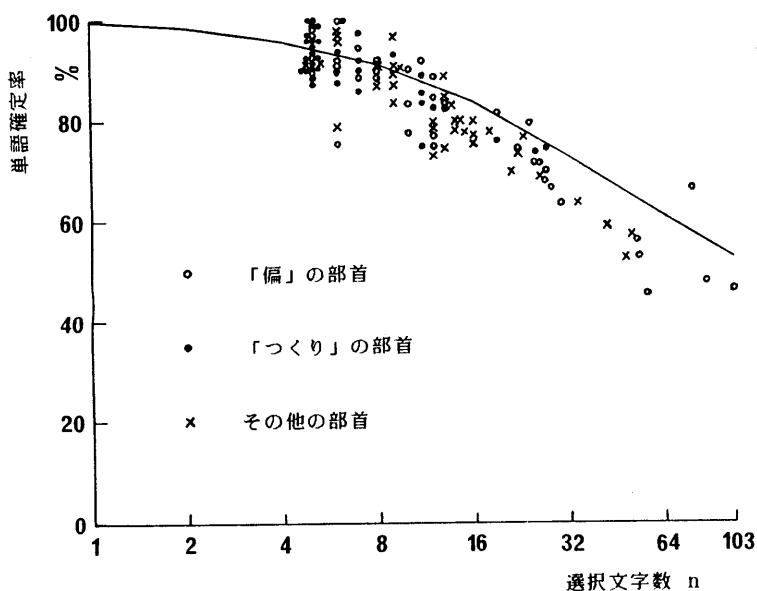


図2 部首共通文字集合の単語確定率

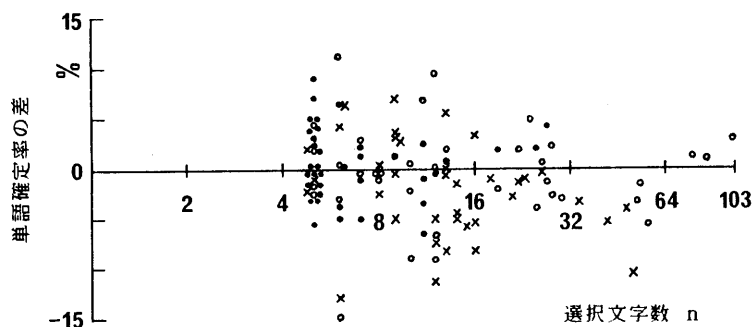


図3 ランダムに候補文字種を選択した時との単語確定率の差

前述の単語辞書について、式(4)から求めた、選択文字数nに対する確定文字率を図1に示す。これより、何らかの方法によって、各文字に対する候補文字種をランダムに8個に絞ることができれば確定文字率は約15%に、さらに4個にまで絞ることができれば約33%になることがわかる。当然のことながら、確定文字率は選択文字数の増加とともに急激に低下していく。

3. 部首共通文字集合の性質

3.1 部首共通文字集合の抽出

部首などの漢字の部分構造の情報をもとに、日本語単語辞書を用いることによって、どのくらいの単語が推定でき、文字が決定できるかを検討するために、まず部首を定義し、共通の部首をもつ文字集合を抽出する。

抽出は人手で行い、共通の部首をもつ漢字が5文字以上存在するもののみを対象とした。なお、漢字の特徴として、「偏」や「つくり」、「垂れ」、「構え」、「によう」に属する部首の定義は比較的明確であるが、「冠」と「脚」に属する部首の概念はあまり明確ではない。つまり、これらの部首では、例えば「くさかんむり」や「たけかんむり」、「れんが」、「したごころ」などのように明確なものもあるが、そうでないものもかなり多い。したがって、ここでは「つくり」に対しては漢字の上から見て、また「脚」に対しては下から見て、同じ部分構造をもつものを部首として定義した。

抽出した部首の種類と、それを部首としてもつ漢字の個数を表2に示す。これより、「偏」に属する部首は37種類、「つくり」に属する部首は38種類、さらに「冠」や「垂れ」などの、その他に属する部首が37種類となる。なお、これらの部首を含む漢字の個数は、延べ1722文字である。

3.2 部首共通文字集合の単語確定率

同じ種類の部首をもつ漢字の集合を部首共通文字集合と呼ぶことにする。表2に掲げた全ての部首

について、部首共通文字集合の単語確定率を求めたものを図2に示す。これは、例えば「さんずいへん」の部首共通文字集合では、「偏」を共有する103個の漢字の各々に対して、同じ集合に属す103個の漢字が候補文字種として選択された時の、この文字集合の単語確定率を求めたものである。したがって、これは部首などの漢字の部分構造が明らかになった時の、つまり漢字の部首情報からの単語確定率に該当する。なお実線は、全ての文字種について、全字種のなかからランダムに候補文字種を選択した時の単語確定率を式(3)から求めたものである。

これより、部首の種類によってばらつきはあるものの、選択文字数が少ない場合には、両者はかなりよく重なることがわかる。しかし、選択文字数が増加するにしたがって、部首共通文字集合の単語確定率の方が小さくなる傾向にある。これは、選択文字数が比較的小さくても単語を決定しやすい単語長の長い単語にひらぐが多く含まれ、それがランダムに選択した時の単語確定率に表れているためと考えられる。

このことをもう少し詳しく検討するために、部首共通文字集合の単語確定率と、それと同じ文字集合の漢字に対してランダムに候補文字種を選択した時の単語確定率との差を求めたものを図3に示す。これより、文字集合によってばらつきはあるものの、全体としては差0を中心として上下にほぼ均等に分布しており、ひらがなを除外して考えると、両者はほぼ一致すると言えよう。つまり、同じ部首の漢字に限っても、単語構成要素としての文字の性質では、他の文字と変わらないことがわかる。

しかし、部首の種類によっては、この差が大きくなるものもある。例えば、「偏」に属す部首では、部首共通文字集合の単語確定率の方が8%以上大きくなるものに、「やまへん」や「かいへん」がある。ちなみに、「やまへん」の文字集合では、それを含む単語数が32個と少ないものの、部首共通文字集合の単語確定率は100%になる。一方、部首共通文字集合の単語確定率の方が8%以上小さくなるものには、「ふねへん」や「にちへん」、「こめへん」がある。特に、「ふねへん」の文字集合では、「般」を別にすると、他の「艦、航、船、艇、舶」はよく似た意味をもち、同じような使われ方をするため、単語確定率は15%も小さくなる。

また、「つくり」に属す部首では、部首共通文字集合の単語確定率の方が8%以上大きくなるものは、わずかに「謁、喝、渴、褻、掲」の漢字に含まれる部首の1個だけで、8%以上小さくなるものはない。さらに、これ以外のその他に属す部首では、8%以上大きくなる部首はなく、逆に、8%以上小さくなるものに「ひとがしら」や「なべぶた」、「あめかんむり」、さらに「宮、栄、覚、学、蛭、労」の漢字に含まれる

部首の4個がある。特に、後ろの二つは10%以上も小さくなっている。

このように、漢字の部首などの部分構造に着目して候補文字種を限定しても、部首によってばらつきはあるものの、全体としては、ランダムに候補文字種を選択した時と、単語としての文字の性質にはあまり違いないことがわかる。つまり、単語としての文字の決定において、漢字の部首情報もそれほど特別の意味をもっているとは言えない。

4. 偏・つくり文字集合の性質

漢字の部首のうち、「偏」と「つくり」の概念は比較的明確であり、かつそれらは相補的な関係にある。つまり、「偏」をもつ漢字には、必ず「つくり」に相当する部首が存在する。一方、「冠」と「脚」もまた相補的な関係にあるが、そこに含まれる部首の概念はあまり明確ではない。すなわち、常用漢字の文字集合を対象に、部首を定義して、「冠」と「脚」に相当する部分構造を矛盾なく抽出することは極めて難しい。また、「垂れ」、「構え」、「にょう」には、相補的な部首の考えはなく、部首以外の他の部分構造は漢字によってほとんど全て異なったものになっている。

したがって、ここでは、「偏」と「つくり」を含む漢字だけの文字集合を対象として、部首などの漢字の部分構造の情報から、どのくらいの単語が決定され、その漢字が推定できるかについて検討する。

4.1 偏・つくり文字集合の抽出

常用漢字1945文字のうち、「偏」と「つくり」を含む漢字は1045文字に及ぶ。ただし、これらの漢字は全て左右に分割でき、かつ「偏」や「つくり」の概念の明確な部首に限定した。したがって、たとえ左右に分割できても、部首の明確でないもの、例えば「八、兆、非、卯、以、川、州、多」などの漢字は除外した。また、「偏」の概念が明確でも、パターンとして左右に分離しにくいと考えられる、「死」も除外した。

こうして抽出した偏・つくり文字集合において、異なる「偏」の種類は196個であった。その主なものは表2に示した通りである。したがって、各部首に含まれる漢字の個数は平均5.3個となる。同じ部首をもつ漢字が最も多く存在するものは「さんずいへん」で、103個にも及ぶ。ついで、「にんべん」の86個、「てへん」の78個の順となる。「偏」は特定のものに集中しており、上位の10部首だけでも全体の半分を越え、かつ上位の19個を取れば漢字の個数は700個に達する。ちなみに、その部首の漢字がただ一つしか存在しないものは126個であった。

一方、異なる「つくり」の種類は、表2に掲げた部

首を主なものとして、合計520個であった。したがって、各部首に含まれる漢字の個数は平均2.0個になる。最も多くの漢字が存在するのは、「りっとう」の27個である。ついで、「ぼくずくり」の25個、「おおがい」の19個の順となっている。「つくり」は極めて種類が多く、しかもばらついた分布になっている。ちなみに、上位の98個で漢字の個数がやっと500個に達し、その合計が700個になるには上位195個の部首が必要である。なお、ただ一つの漢字しか含まない部首は305個であった。

4.2 偏・つくり文字集合の単語確定率

偏・つくり文字集合を成す各漢字に対して、それと同じ部首をもつ全ての漢字が候補文字種として選択される場合の、この文字集合の単語確定率を求める。

まず、「偏」に着目して、漢字の「偏」の情報が正しくえられたとすると、例えば「さんずいへん」の漢字に対しては同じ部首をもつ103個の漢字が、また「にんべん」の漢字には86個の同じ部首をもつ漢字が、各々候補文字種として選択されることになる。したがって、「加」のようにただ1個の漢字しかない部首では、他に候補文字種が存在しないため、この漢字を含む単語は全て決定されることになる。このようにして、偏・つくり文字集合に含まれる全ての漢字に対する候補文字種を設定することにより、この文字集合の単語確定率を求めることができる。前述の単語辞書において、「偏」に着目した偏・つくり文字集合の単語確定率は72.43%であった。すなわち、漢字の

「偏」が何であるかがわかることにより、それらを含む単語の3/4近くが決定できることになる。

同様に、漢字の「つくり」の情報が正しくえられたとすると、この文字集合の単語確定率は94.57%になった。つまり、漢字の「つくり」がわかると、それを含む単語の約95%を決定できることになる。単語辞書の他に、さらに文節や文脈の情報をも利用すると、ほとんど全ての単語を決定できる可能性がある。

4.3 形状類似文字集合の単語確定率

つぎに、文字全体の形に着目して、形状類似文字を候補文字種として選択した時の、単語確定率について検討する。形状類似文字集合を機械的に抽出する確固たる理論はないが、特徴抽出が単純明快で、かつ手書き漢字に対して比較的高い認識率のえられているものが文字の形状をよりよく表しているものと考え、ここでは外郭方向寄与度特徴[8]を用いて、形状類似文字集合を求めた。手書き文字パターンは、電総研で作成されたJIS第1水準手書き漢字データベースETL9[9]を使用し、各文字種について200サンプルでの平均外郭方向寄与度特徴を求め、L₁距離を尺度として類似文字集合を抽出した。

いま、偏・つくり文字集合の「偏」に着目すると、例えば「さんずいへん」の103個の漢字に対しては各々103個の候補文字種を、「にんべん」の86個の漢字は各々86個の候補文字種をもつため、この集合の漢字全体に対する平均候補文字数は38.93個となる。そこで、偏・つくり文字集合の各漢字に対して各々39個の形状類似文字を選択し、その単語確定率を求めると65.18%であった。よって、「偏」に属す部首に着目した前述の単語確定率の方が、漢字の形状に基づいて同数の類似文字種を選択するよりも、約7.3%大きくなり、より多くの単語を決定できることがわかる。

同様に、「つくり」に着目すると、偏・つくり文字集合の漢字に対する平均候補文字数は4.70個になる。そこで、各漢字に対して4~5個の形状類似文字種を選択して単語確定率を求め、それから内挿すると、この単語確定率は約94.0%になった。つまり、「つくり」についても、全体形状に基づいて類似文字種を選択するより、同じ部首の漢字を選択する方がわずかではあるが単語確定率は大きくなるのがわかる。

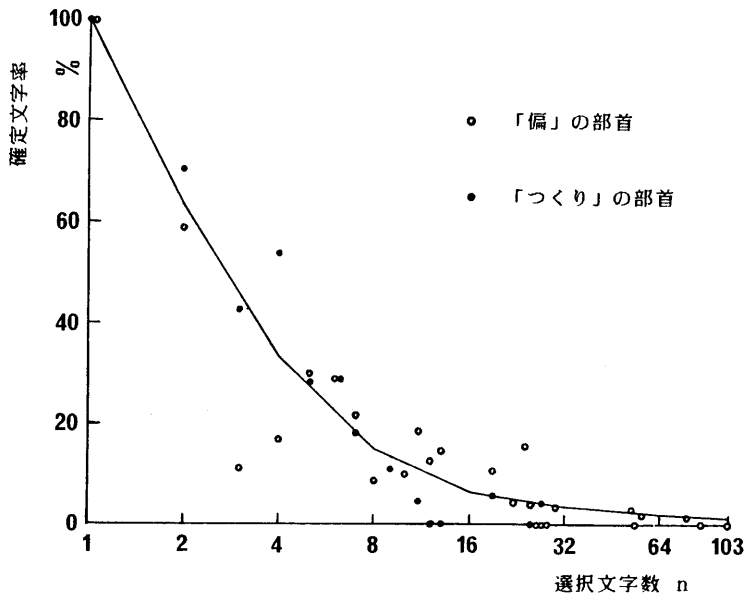


図4 部首共通文字集合の確定文字率

4. 4 偏・つくり文字集合の確定文字率

まず、「偏」または「つくり」に属す部首共通文字集合の確定文字率について求めたものを図4に示す。これは、同じ部首の全ての漢字を候補文字種としても単語を完全に決定できるものを求め、「偏」あるいは「つくり」のうちで選択文字数を同じにする漢字群を改めて一つの集合とし、確定文字率を算出したものである。また実線は、全ての文字種を対象として、全字種のなかからランダムに候補文字種を選択した時の確定文字率を式(4)から求めたものである。

部首共通文字集合の確定文字率は、算出の対象となる文字数が限られるため、ばらつきが大きくなるものの、ランダムに候補文字種を選択した時のそれによく重なっていることがわかる。つまり、確定文字率においても、同じ部首の文字集合に、単語としての特別な性質の違いはないことがわかる。

つぎに、偏・つくり文字集合の確定文字率について検討する。単語確定率の場合と同様に、まず「偏」に着目し、各漢字に対して「偏」に属す同じ部首の漢字が候補文字種として選択されるとすると、この文字集合の確定文字率は19.23%であった。この結果は漢字の「偏」が何であるかがわかることにより、その漢字を含む全ての単語が推定でき、どの単語からでもその漢字を一意に決定できるものが約2割に達することを示している。なお、この場合の平均選択文字数に相当する39個の文字を全字種からランダムに選択した時の確定文字率は、式(4)より求めると、わずかに2.85%である。

つぎに、「つくり」に着目して、偏・つくり文字集合の確定文字率を求めると58.28%であった。漢字の「つくり」がわかると、6割近くの漢字がそれを含むどの単語からでも一意に決定できる性質をもっていることがわかる。なお、平均選択文字数に相当する4.7個の候補文字種をランダムに選択した時の確定文字率は約28.0%である。このように、一定個数の候補文字種をランダムに選択するよりも、「偏」または「つくり」に基づいて候補文字種を選択する方が確定文字率ははるかに大きくなる。

5. むすび

本報告では、部首などの漢字のもつ部分構造がわかると、単語辞書を使用することによって、どのくらいの単語が決定でき、文字が一意に推定できるかを単語確定率と確定文字率を評価尺度として定量化した。その結果、つぎの事項が明らかとなった。

(1) 112種類の部首を定義して抽出した部首共通文字集合の単語確定率は、同じ個数の候補文字種を全ての文字種のなかからランダムに選択した時のそれとほぼ同じになる。ただし、部首の種類によっては、この差が大きくなるものもある。

(2) 「偏」と「つくり」から構成されている漢字の「偏」が何であるかがわかると、それらの漢字を含む単語の約72.4%が決定できる。同様に、漢字の「つくり」がわかると、その単語確定率は約94.6%になり、極めて高い値を示す。

(3) これらの単語確定率は、漢字の全体形状に着目して、これと同数の類似文字種を一律に選択するよりも大きくなる。特に、「偏」に属す部首についての単語確定率は約7.3%も大きくなる。

(4) 「偏」および「つくり」に属す部首を対象とした、部首共通文字集合の確定文字率も、単語確定率と同様に、同数の候補文字種をランダムに選択した時のそれとほぼ同じになる。

(5) 「偏」と「つくり」から成る漢字の「偏」がわかると、その確定文字率は約19.2%になる。同様に、漢字の「つくり」がわかると、その確定文字率は約58.3%に達する。すなわち、6割近くの漢字がそれを含むどの単語からでも一意に決定できる性質をもっている。

このように、漢字の「偏」が何であるかが明らかになると、単語辞書を用いることによって、約72%の単語が決定でき、2割近くの漢字については、それを含む全ての単語が決定できる。一方、「つくり」がわかると、約95%にも及ぶ単語が決定でき、6割近くの漢字については、それを含む全ての単語を決定することができ、どの単語からでもその漢字を一意に決定できることになる。

参考文献

- [1] 王, 真田, 手塚: "手書き文字認識について", 信学技報, PRL84-81, pp.35-44(昭60-02).
- [2] 津雲, 浅井: "文字認識技術の最近の動向", 信学技報, IE88-5, pp.31-38(昭63-04).
- [3] 小川, 手塚: "漢字の階層表現とその認識", 信学論(D), J57-D, 12, pp.700-707(昭49-12).
- [4] 吉田: "ストローク抽出法による手書き漢字認識システム", 信学全大, 1543(昭49-03).
- [5] 安居院, 中嶋, 長橋: "部分パターンの位置関係を利用した手書き漢字の表現法", 信学論(D), J60-D, 12, pp.1109-1116(昭52-12).
- [6] 金田一, 池田編: "学研国語大辞典", 学習研究社.
- [7] 梅田: "単語辞書を用いた文字認識における文字の確定能力", 信学論(D-II), J72-D-II, 1, pp.22-31(平1-01).
- [8] 萩田, 内藤, 増田: "外郭方向寄与度特徴による手書き漢字の識別", 信学論(D), J66-D, 10, pp.1185-1192(昭58-10).
- [9] 齊藤, 山田, 山本: "JIS第1水準手書き漢字データベースETL9とその解析", 信学論(D), J68-D, 4, pp.757-764(昭60-04).