

# コスト最小法を用いたことば遊び ～数字語呂合わせの自動生成システム～

青木賢太郎

慶應義塾大学大学院 政策・メディア研究科

**Abstract:** 形態素解析に用いられるコスト最小法で計算されるパスコストは文の文としての妥当性を表していると言える。そこで、本研究では各数字に割り当てられた単純な語呂から、長い数字列の持つ文としての妥当性をパスコストによってはかり、語呂合わせを自動生成するシステムを作成した。また、形態素解析器の参照する辞書の中身をあらかじめ数字列に置き換える事によって、入力された数字列を直接、形態素解析可能にし高速化を図った。

## 1 はじめに

語呂合わせの歴史というのは非常に古く、遡ればことばの読みそのものが意味や力を持つという言霊の時代にまで行きつく。そこまでいかないにしても、江戸時代の人々は好んで語呂を合わせ、歌や落語などの中でその巧さを競っていた。現代では、語呂合わせは主に何かを覚えるために使われている。百人一首、元素周期表、中国王朝などもあるが、我々の日常で一番良く使われるのが、歴史の年号、電話番号、円周率などの数字語呂合わせである。覚えるためだけでなく、覚えさせるためにも数字語呂合わせは使われている。例えば、企業は自社の電話番号をできるだけ意味の通った語呂合わせで申請し、祝日や記念日がそれに関連した語呂になっていることも多い。欧米の電話番号などが、アルファベットを割り当てたりして、なんとか意味をつけて覚えやすくしようとしているのに対し、日本では数字の読みの集合から別の語を連想することによって記憶するという方法が取られている。

数字語呂合わせは日本だけではなく中国語圏などでも行われている。例えば広東語で4は「セイ」と発音し、「死」と同じ音となり一般には、不吉な数字として忌み嫌われる。実は、これと同じ現象が古代日本で起

きた結果を、我々は普段口にしていて、我々は数字を「イチ、ニ、サン、よん、ゴ、ロク、なな、ハチ、キュウ、ジュウ」と数えることがあるが、これは呉音をベースとした読みで、「シ、ク」は「死、苦」の表現につながることから、「シ、シチ、ク」を和語読みの「よん、なな」と漢音の「キュウ」に置き換えたためである。[1] (表1 参照)

数字の語呂は現在でも縁起を担ぐために使われるが、やはり数字列を覚えるために使用されるのが身近である。我々が語呂合わせをする際、通常図1のようなプロセスを経る。

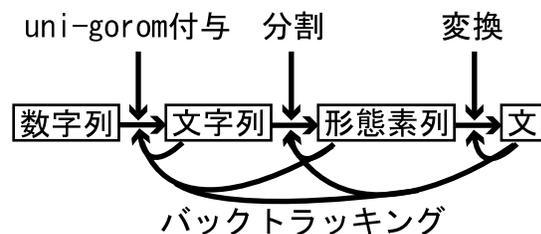


図1: 語呂合わせのモデル

まず、数字列を各数字が持つuni-gramの語呂を用いて文字列に変換する。ここで、便宜上または語呂の良さからuni-gramの語呂のことを**uni-gorom**と呼ぶ。

続いてその文字列を形態素に分ち書きし、文として意味のあるものに変換するわ

けだが、その解析過程において何度もバックトラッキングが発生しているのが分かる。また、各過程において生成される組み合わせの総数は爆発的であり、その全ての経路を把握することは人間では不可能である。そこで本研究では、コスト最小法を用いた形態素解析によって、数字語呂合わせを実現するシステムを作成した。

## 2 関連研究

前節で述べたように、数字の語呂合わせは負担が大きく、またニーズも多いため、これを解決する方法はいくつか提案されている。しかしその何れも、自動生成というよりは、人間によるバックトラッキングの補助といった色合いの強いものである。

### 2.1 半自動システム

半自動のシステムには、WWW 上で CGI として動作するものとスタンドアローンのアプリケーションとして動作するものがある。前者は主に携帯電話のインターネット接続サービスから、電話番号の語呂合わせをするものであり、後者は「ごろさく」[2]というソフトでこちらも電話番号を意識して12桁までの入力となっている。「ごろさく」は語呂辞書に登録された3000語以上の単語の組み合わせを効率良く提示できるが、CGIのものと同じく区切りはユーザが指定するようになっている。つまり多くの文節候補の中から、組み合わせを自分で選ぶという、最後の部分が手付かずである。

### 2.2 自動システム（プロトタイプ）

この他に、筆者による本システムのプロトタイプがある。本研究は言わばこれの改良手法にあたる。このシステムはコスト最小法に基づいて、数字列を解析する点では変わっていないが uni-gorom 付与の過程（図

1)をそのままなぞっているため、最初に多くの文字列を生成していた。そのため、

- 中間ファイルが非常に大きい
- 時間がかかる（4桁で1分、5桁は待てない程）
- uni-gorom の語呂を分割してしまう（279→にしちく→西地区）

のような問題を抱えていた。

## 3 システムの特徴

本システムの要素技術は形態素解析に用いられるコスト最小法という技術であり、そのエンジンとして既存の日本語形態素解析器である ChaSen（茶筌）[3]を用いた。茶筌が参照する辞書の見出し語を予め数字列で置き換えておくことによって、直接数字列を解析することが可能になるわけである。そのため、数字列の桁数を気にせず高速な生成ができる。

### 3.1 語呂

辞書の生成に当たり、語呂の選定は非常に重要である。数字の語呂合わせにおいて、語呂が付与されるのは正確には数詞という品詞に分類される語である。その数詞は以下の様に分類される。[4]

- 基数詞 – 事物の数を表す語
- 序数詞 – 順序を表す語
- 倍数詞・分数詞 – 英語の double や half など

その他に集合数詞、配分数詞、不定数詞などがあるが、日本語では基数詞、序数詞が主であり、読みに対する語呂もこれらに付与されるのでここでは詳しく述べない。



(品詞 (名詞 一般)) ((見出し語 (異例 2720)) (読み イレイ) (発音 イレイ))

↓↓↓↓↓

(品詞 (名詞 一般)) ((見出し語 (1 0 2720)) (読み 異例) (発音 イレイ))

図 4: システム辞書と変更後

以上の変更を行った結果、元の辞書の項目が23万あまりであったのに対し、語呂合わせ用の辞書では4万5千あまりと5分の1程になった。これは数字からは変換できない仮名が存在するためである。(表2参照)だが、「ごろさく」の3000語と比較すると15倍もの開きがあり、計算機が組み合わせを求める能力との差を歴然と物語る。

## 4 実行結果

実行結果を図5に示す。これは茶釜の解析結果をグラフィカルに表示する Visual Morphs[5]を使った例である。パスコストを見ながら自分の好きな経路を選べるようになっている。

## 5 発展

発展としては

- uni-gorom を選べる
- 概念ベースと組み合わせる

などが考えられる。前者では自分の馴染みの uni-gorom で解析することによって語呂を覚えやすくする効果が期待できる。

また、後者では予めキーワードを入力し、概念ベースを用いて類似度を計算するなどしてユーザの趣向にあった語呂を優先的に出力できるものと考えられる。(例 4 1 2 9 - ブティックには「良い肉」よりも「良い服」など)

## 6 おわりに

本稿では、形態素解析に用いられるコスト最小法のアルゴリズムによって数字の語呂合わせを自動的に生成するシステムについて述べた。人間ではとても思いつかない量の語呂を生成し、それをパスコストという文の妥当性を示す指標で順位付ける事が可能になり、ユーザの負担が大幅に減少した。

今後はユーザの嗜好によるカスタマイズについて、また中国語版の辞書作成などについて検討したい。

## 謝辞

茶釜の Tips を教えてくださった石崎研究室博士課程の金子拓也さんに感謝致します。お陰様で道が開けました。またシステム全般において手伝ってくれた花川賢司君に感謝致します。

## 参考文献

- [1] 三省堂編修所 編:ことばの知識百科, 三省堂,1995
- [2] todome:ごろさく, <http://hp.vector.co.jp/authors/VA001042/>
- [3] 松本裕治, 北内啓, 山下達雄, 平野善隆, 松田寛, 高岡一馬, 浅原 正幸: 日本語形態素解析システム『茶釜』version 2.2.1 使用説明書, Nara Institute of Science and Technology, 2000
- [4] 小川芳男, 林大: 日本語教育事典, 大修館書店, 1982
- [5] 松田 寛, 松本 裕治: 品詞タグ付きコーパス作成支援 GUI ツール Visual Morphs, 情報処理学会研究報告 2000-NL-137, p.98, June, 2000

音読み	0	1	2	3	4	5	6	7	8	9	10
呉音		いち	に	さん	し	ご	ろく	しち	はち	く	じゅう
漢音	れい	いつ	じ	さん	し	ご	りく	しつ		きゅう	
慣用音											じっ
訓読み											
		ひと	ふた	み	よ	いつ	む	なな	や	ここの	とお
		ひとつ	ふたつ	みつ	よつ	いつつ	むつ	ななつ	やつ	ここのつ	
				みつつ	よつつ		むつつ		やつつ		
	ゼロ	はじめ			よん		むい	なの	よう		

表 1: 数詞の日本語読み

ん	わ	ら	や	ま	は	な	た	さ	か	あ	*
ん	0	ら	8	0	8	7	た	3	か	2	a
*	ゐ	2	1	3	1	2	7	4	9	1	i
*	*	る	ゆ	6	2	ぬ	2	4	9	5	u
*	ゑ	0	8	め	へ	ね	10	7	け	8	e
*	0	6	4	も	ほ	0	10	3	5	0	o

表 2: 50 音対応表

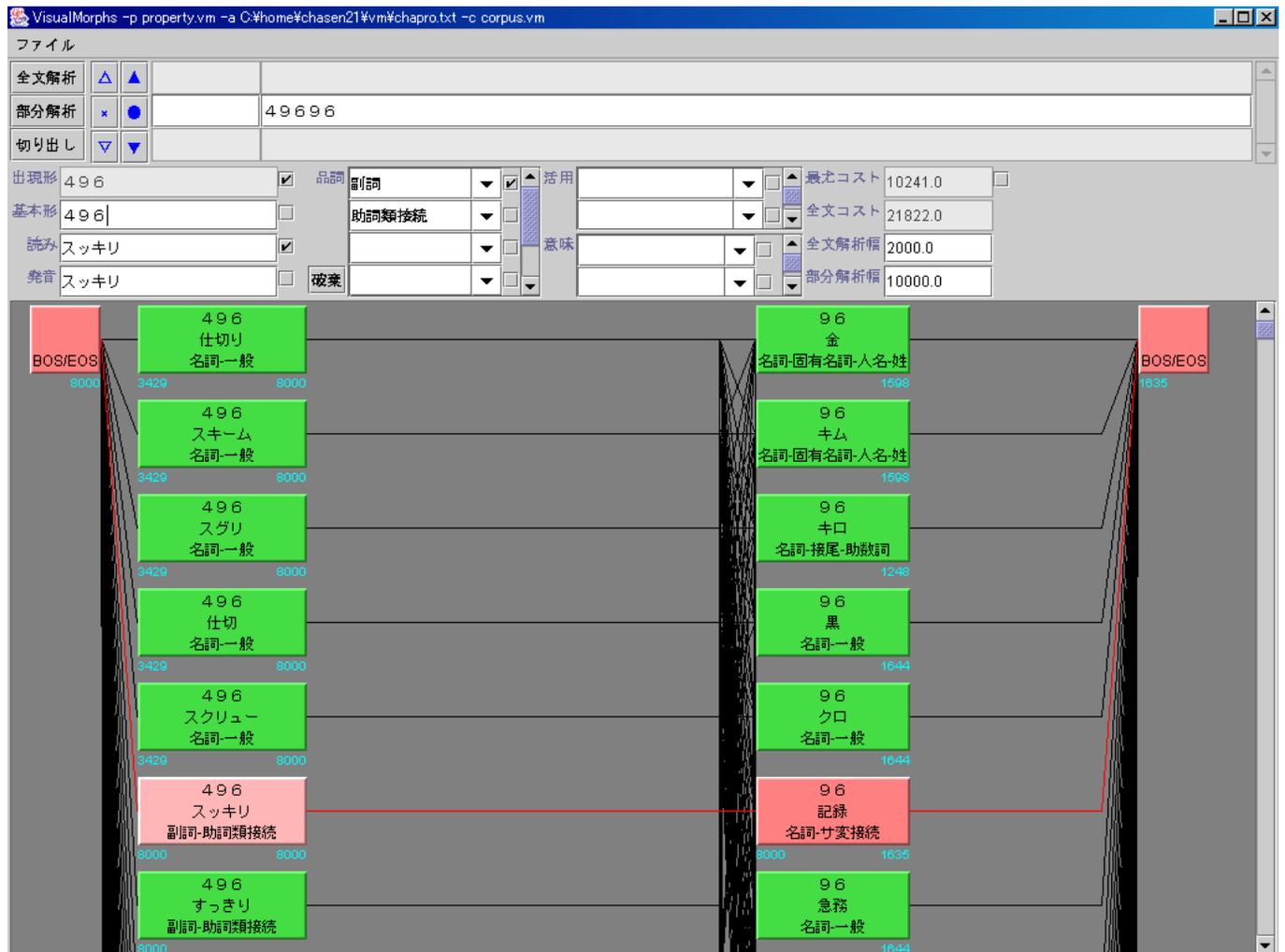


図 5: Visual Morphs を使った実行例