

## 日本語の形態素解析について

首藤 公昭  
(福岡大学・工学部)

## 1. まえがき

表現形式や意味内容の多様性は、自然言語の特徴であり、機械処理においてこの問題を避けて通ることはできない。筆者らは、日常の口語文章体文(技術論文等)をこのまゝ扱う日本語解析システムを開発すべく、基礎となる言語情報の整理を行ってきた。

量的にほう大となる体言、用言(及びこれらに相当する複合表現)に関する情報と文のわく組みに関する情報とを分離し、前者は、処理対象の分野ごとにリフレッシュして利用する一方、後者については、網ろ的を整理を行って、分野に独立で強力な資源を与えておくことを考えている。すなわち、分野は限定されるが、取扱う文のわく組みには大きな自由度を持つたシステムの開発をねらいとしている。本稿では、おもに後者に関する総合的な報告を行う。

日常の日本語文には、文法カテゴリーの判然としない表現や慣用句的表現がかなり頻繁に用いられている。また、日本語は、改米語と異なり、語の屈折より膠着によって構文情報が与えられるという特徴を持つが、日常の文では、膠着の仕方はかなり複雑で、従来、以上をまとめて体系的に整理する研究は見られない。我々は、意味処理への見直しを重視しつつ、文のわく組みを構成する表現を網ろし、文の直接構成単位と見なし得る局所的表現列の構造を明らかにした。この局所的表現列は、「文節」の考えを意味のうえから拡張したものであり、我々は、「E-文節」と呼んでいる。改米語における語の

屈折を取扱う段階との類似性から、E-文節の構造を解析する段階を(広義の)「形態素解析」と呼ぶこともできよう。

なお、日本語にも膠着(以後、接続という)の仕方と密接に関連した屈折、すなわち「活用」があり、我々は、これに関する機械処理向きな整理を行い、ほとんど例外なく処理できる見通しを得ている。

## 2. 要素表現の設定

意味機能の点で単語的に扱うのが望ましいと思われる複合表現も含めて、分野に余り関係なく使用され、量もさほどほう大とはならない表現を網ろ的に収集した。資料は、高校の教科書、技術論文の口語文章体、約13,000文である。

## 2.1 付属語的表現

日本語は膠着言語、あるいは後置詞言語であり、文のわく組みは、後置詞、すなわち付属語による。すなわちと云われる。しかし、意味の点から日常の文を見れば、付属語と類似の機能を有する複合表現が多数見出される。例えば、

「それについても考えておかなければならないのではなかろうか。」

という文でいわゆる付属語を文頭側から抽出すれば、「に、て、も、ぞ、ない、は、ない、の、で、は、ない、う、か」となるが、この様を解析を行うのは、文の意味を把握する上で得策でない。

そこで、人間が理解する際に一まじりに認識すると思われる「について、も、ておく、なければならぬ、ので

はなかるうか'等を付属語的な要素表現とみなす。この種の表現が日常、どれくらい使われるかが問題であるが我々は、前記の資料から約950種を得、その後の調査からかなり網らできたのではないかと考えている。これらの表現のうち、'についで'、'における'など格助詞、副助詞的なものを「関係表現」、'なければならぬ'、'のではなからうか'など助動詞、終助詞的なものを「助述表現」と呼んでいる。

広範な日本語文を扱うシステムを考える際、前者は、日本語の表層格、後者は、モーダリティ、アスペクト等を与える重要な表現群である。

2.2 接尾語的表現

日常の文には、多数の接尾語的表現が用いられるが、これらのうち、分野に余り関係なく、高い頻度で用いられるもの約120種を抽出した。例えば、'調査するかどうかは問題なのである'、という文からは、'する'、'かどうか'、'なのである'を抽出した。その他の接尾語については、接頭語、助動詞等と同様、分野を定めて整理することを考えている。

2.3 自立語的表現

接続詞、連体詞、副詞的表現の抽出を行った。ここでも'このような'、'しかしながら'、'従って'、'ある程度'などの複合表現を含め、計約400種を定めたが、必ずしも十分とは言えないため、現在、増補中である。

3. 文の局所的表現列(E-文節)の構造

3.1 要素表現の分類

2.で定めた要素表現、特に付属語、接尾語的表現の接続機能を詳細に調べ、これに基づいて分類を行って約140個の文法カテゴリーを定めた。

これは、学校文法での文法カテゴリーを接続機能によってリファインしたもので、同綴り異機能(異義)の要素表現は、機能ごとに別表現とみなした分類となっており、接続の解析に下って要素表現の多義選択が可能となる様、配慮している。

分類の概要を表1に示す。(表現列については図1を参照のこと)

表1 要素表現の分類

大分類	表現数	文法カテゴリー (code)	
付属語的表現	R <sub>NP1</sub>	155 R01 R09, R10~R13, R17~R19, R20~R27	
	R <sub>NP2</sub>	65 R30~R39, R40~R46, R50~R59, R60~R63	
	R <sub>NP3</sub>	16 R69	
	R <sub>NN1</sub>	1 R70	
	R <sub>NN2</sub>	119 R71	
	R <sub>NN3</sub>	4 R72	
	R <sub>PN</sub>	38 R80	
	R <sub>PP1</sub>	150 R90~R93	
	R <sub>PP2</sub>	1 R94	
	R <sub>PP3</sub>	39 R95~R97	
	助述表現	*A <sub>NP1</sub>	41 A11, A15, A18, A1D, A14, A17
		*A <sub>NP2</sub>	4 A24, A29
		*A <sub>NP3</sub>	3 A30
		*A <sub>PP1</sub>	299 A40, A41, A43~A49, A4A~A4D
*A <sub>PP2</sub>		13 A50	
*A <sub>PP3</sub>		4 A60	
*A <sub>PP4</sub>		1 A70	
(S <sub>nn1</sub> )		290 S50~S59, S60~S69, S70~S74)	
S <sub>nn2</sub>		6 S80	
S <sub>nn3</sub>		4 S90, S91	
接尾語的表現	S <sub>pn</sub>	92 S40~S49	
	*S <sub>np1</sub>	12 S21, S29, S2A, S2D	
	*S <sub>np2</sub>	2 S35, S34	
	*S <sub>pp</sub>	2 S11, S1D	
	自立語的表現	C <sub>1</sub>	36 C10
		C <sub>2</sub>	87 C20
		D <sub>1</sub>	185 D10
		D <sub>2</sub>	80 D20
		T	164 T00
		(M <sub>1</sub> )	- M11~M15)
(M <sub>2</sub> )		- M20)	
用言	Y	- Y00~Y09, Y0A, Y0B)	

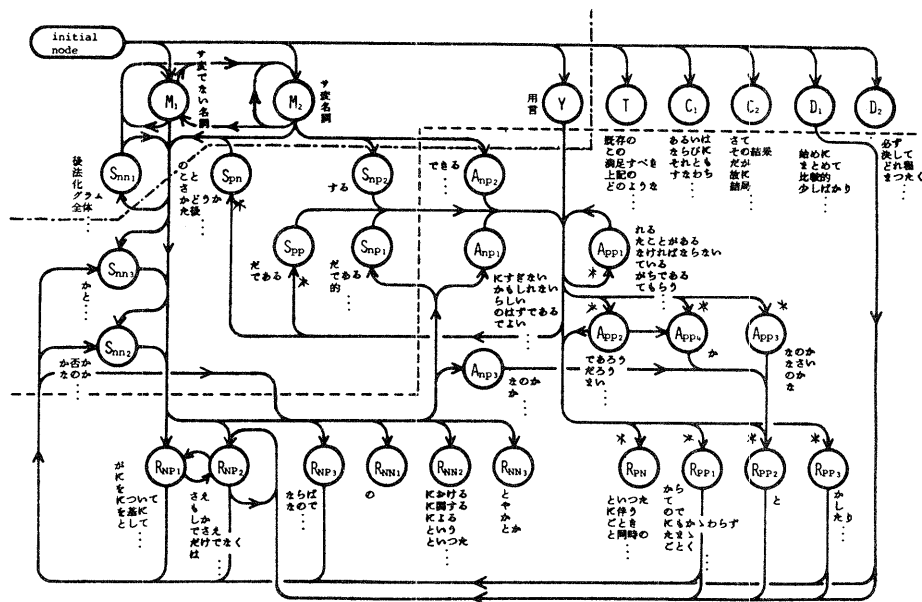


図1 接続ルールの概要

れば、不変化部 'KIS' を語幹とすることができ

(2) 使役、可能動詞の語尾を切離し、独立した助述表現とみなすことができる。

例えば、使役動詞 '降らす' は、Fu

### 3.2 接続ルール

3.1で定めた各カテゴリーに対して、どのカテゴリーに接続可能かを規定するルールを与えた。ルール数は、被接続、接続カテゴリーの対、約3,600個である。ルールの詳細は略すが、その概要を表1の大分類のカテゴリー32種を用いて、遷移図の形で図1に示す。

ノードは、大分類でのカテゴリー、矢印は起算のノードの表現に着算のノードの表現が文末側から接続できることを表す。また、ノードの遷移の過程で、破線を上部の領域から下部の領域へ最後に横切る時が、E-文節の自立部と付属部の境界に大よそ対応している。

### 3.3 ローマ字表記

2.で抽出した各種表現は、現在、ローマ字表記によっている。

ローマ字表記には、次の利点がある。(1) 語幹、活用語尾の分離が音素レベルで行え、活用のルール化が容易になる。

例えば、変動詞 '期する' の連用形は '期し' となるが、ローマ字表記によ

RASU' となり、語幹 'FUR' に使役の助述表現 'ASU' が接続したものとす。

また、可能動詞 '行ける' は、'行く' の語幹 'IK' に助述表現 'ERU' が接続したものとす。

(3) 英字入出力機器が使える。

(4) 仮名との相互変換が容易である。

(5) 表記としての国際性。

筆者らは、特に、(1)、(2)を重視し、各行で子音部文字が変化しない表記法を採用した。これを表2に示す。

表2 採用したローマ字表記

A	I	U	E	O				
ア	イ	ウ	エ	オ				
KA	KI	KU	KE	KO	KYA	KYU	KYO	
カ	キ	ク	ケ	コ	キヤ	キユ	キョ	
SA	SI	SU	SE	SO	SHA	SHU	SHE	SHO
サ	シ	ス	セ	ソ	シャ	シュ	シェ	ショ
TA	TI	TU	TE	TO	CHA	CHU	CHE	CHO
タ	チ	ツ	テ	ト	チャ	チュ	チェ	チョ
NA	NI	NU	NE	NO	NYA	NYU	NYO	
ナ	ニ	ヌ	ネ	ノ	ニヤ	ニユ	ニョ	
HA	HI	HU	HE	HO	HYA	HYU	HYO	
ハ	ヒ	フ	ヘ	ホ	ヒヤ	ヒユ	ヒョ	
MA	MI	MU	ME	MO	MYA	MYU	MYO	
マ	ミ	ム	メ	モ	ミヤ	ミユ	ミョ	
YA	YU	YO						
ヤ	ユ	ヨ						
RA	RI	RU	RE	RO	RYA	RYU	RYO	
ラ	リ	ル	レ	ロ	リヤ	リユ	リョ	
WA					WO			
ワ					ョ			
Q								
ン								
GA	GI	GU	GE	GO	GYA	GYU	GYO	
ガ	ギ	グ	ゲ	ゴ	ギヤ	ギユ	ギョ	
ZA	ZI	ZU	ZE	ZO	JA	JU	JE	JO
ザ	ジ	ズ	ゼ	ゾ	ジャ	ジュ	ジェ	ジョ
DA	DI	DU	DE	DO				
ダ	ヂ	ヅ	デ	ド				
BA	BI	BU	BE	BO	BYA	BYU	BYO	
バ	ビ	ブ	ベ	ボ	ビヤ	ビユ	ビョ	
PA	PI	PU	PE	PO	PYA	PYU	PYO	
パ	ピ	プ	ペ	ポ	ピヤ	ピユ	ピョ	

これは、ハボン式と日本式を組合せたものである。ハボン式単独の場合、例えば、動詞‘立つ’の未然、連用、終止形は、それぞれ‘TATA’、‘TACHI’、‘TATSU’となって、語幹部が取りにくい。しかし、日本式によれば、‘TATA’、‘TACHI’、‘TATSU’の不変部‘TAT’を語幹とし、語尾‘A’、‘I’、‘U’を他の動詞、例えば‘起こす’などの語尾変化のパターンと共通にすることができる。

### 3.4 活用の整理

上記の表記法に基づき、活用語尾、活用形、活用型を表すのようまとめた。表中、\*印は、音便変化した1文字を表す。

ここでの整理の特徴を列挙すれば、次の通りである。

- (1) 通常、口語文章体文を対象とするため、ウ音便を除外している(また、現時点では、‘ぞす’、‘ます’等の丁寧表現も考慮していない)が、この範囲ではほとんど例外なく処理できる。
- (2) 音便の整理(後述)に基づき、5段

表4 音便変化する文字

活用型	見出し語の 親1文字	入カ文字列中 の対応の1文字	見出し語の例	入カの例
0	G K R	I	TUG(繰ぐ) KIK(聞く) OKOR(起る)	TUIDA KIITA OKOTTA
1	K T	T	IK(行く) KAT(勝つ)	ITTA KATTA
2	N B M	Q	SIN(死ぬ) YOB(呼ぶ) YOM(読む)	SIQDA YOQDA YOQDA

- 型にI, T, Q型の3種を加えた。  
 (3) ‘言う’などに特定の活用型(7型)を与え、‘IWU’などの変化を避けた。  
 (4) ‘ず’、‘まい’、‘う’などを非活用表現とみなし、活用処理から除外した。

### 3.5 音便の整理

音便変化する語幹の末尾1文字を表4のように整理した。

### 3.6 活用に関する接続条件

表1で\*印を付したカテゴリーの表現は活用するものであり、これらに文末側から何らかの表現が接続する際(図1で\*印の接続)、特定の活用型、および活用形であることが要求される。

3.2で述べた接続ルールには、この種

表3 活用の整理

活用形		活用型										語幹			
活用型	Code	1	2	3	4	5	6	7	8	9	A	B			
動詞型	5段音便I型	0	A	A	0	I	I*	U	U	E	E			ex. KIK(聞く)	
	5段音便T型	1	A	A	0	I	T*	U	U	E	E			ex. OKOR(起る)	
	5段音便Q型	2	A	A	0	I	Q*	U	U	E	E			ex. SIN(死ぬ)	
	5段型	3	A	A	0	I	I	U	U	E	E			ex. KES(消す)	
	1段型	4	e	e	YO	e	e	RU	RU	RE	RO	YO		ex. TOZI(肉じ)	
	サ変型	5	I	E	IYO	I	I	URU	URU	URE	IRO	EYO		S(す)	
	カ変型	6	0	0	OYO	I	I	URU	URU	URE	OI			K(来る)	
W型	7	WA	WA	0	I	T	U	U	E	E			ex. TIGA(違う)		
形容詞型	8	KU		KARO	KU	KAT	I	I	KERE				ex. YO(良い)		
形容動詞型	NA型	9				NI			NA				ex. KIREI(きれい)		
	NO型	A				NI	e		NO				ex. HODO(ほど)		
	E型	B				NI			e				ex. ONAZI(同じ)		
	T型	C	E		ARO			A	A				T(た)		
	D型	D	E	EHA	ARO	E	AT	A	A				D(た)		

の制約は含まれておらず、接続ルール表と切り離して、活用表現に文末側から接続可能な各表現の辞書項目として与える。

この項目は、何活用の何形に接続可能な表すに基づいて記した

もので、表5に示す45種のいずれか一つをコードで与えておく。例えば、助述表現「ているところだ」に対しては、K-43が記され、2型を除く動詞型活用表現の活用形5に接続することが示される。

この条件をチェックするには、被接続の活用表現がどの活用型であるかを知らねばならないが、表1に示した活用表現の小分類は、活用型による分類となっており、コード3桁中、最下位桁は表3に示した活用型コードその

ものである。例えば、用言の動詞、形容詞、形容動詞への細分類は、活用型によって行われていることになる。

#### 4. 辞書

2.で定めた要素表現を収録した機械辞書を作成した。見出し語に続く文法情報の項目は、表1の文法カテゴリコード、および表5の活用接続条件コードの2種である。一部を図2に示す。

表5 活用接続条件

Code	接続可能な活用型、活用形
K-01	(5 B)
K-02	(0-3 B)
K-03	(0-4 B) (7 B)
K-06	(8 B)
K-07	(A B)
K-08	(9-B B)
K-09	(8-B B)
K-11	(0-3 1) (7 1)
K-12	(0-7 1)
K-13	(0-8 1)
K-14	(0-8 1) (D 1)
K-15	(0-8 1) (D 1) (D 2)
K-21	(0-7 2)
K-22	(4 2) (6 2)
K-25	(0-7 3)
K-26	(0-8 3)
K-27	(0-8 3) (D 3)
K-31	(0-7 4)
K-33	(0-7 4) (8 B)
K-34	(0-7 4) (8-B B)
K-36	(0-B 4)
K-41	(0 5) (2 5)
K-42	(0 5) (2 5) (9-B B)
K-43	(0-1 5) (3-7 5)
K-45	(0-1 5) (3-7 5) (8 4)
K-46	(0-1 5) (3-8 5) (D 5)
K-51	(0-7 6)
K-52	(0-7 6) (C-D 6)
K-53	(0-8 6)
K-54	(0-8 6) (9-B B)
K-55	(0-8 6) (9-B B) (C-D 6)
K-56	(0-8 6) (C-D 6)
K-57	(0-8 6) (D 6)
K-61	(0-7 7)
K-62	(0-7 7) (C-D 7)
K-64	(0-8 7) (C-D 7)
K-66	(0-8 7) (9-B B) (C-D 7)
K-67	(0-9 7)
K-68	(0-9 7) (B-D 7)
K-69	(0-9 7) (C-D 7)
K-70	(0-B 7)
K-71	(0-D 7)
K-72	(8-D 7)
K-73	(C-D 7)
K-78	(0-8 8)

0230 00	TOHAKAGIRANA	C-A48	K-55
0231 12	TOHAKAGIRAZU	C-R95	K-55
0232 12	TOHIKAKUSITE	C-R50	
0233 12	TOIUKATATIDE	C-R17	
0234 12	TOIUKEITAIDE	C-R17	
0235 12	TOKAQKEINAKU	C-R09	
0236 12	TOKAQREQSITA	C-R71	
0237 12	TOKAQREQSITE	C-R17	
0238 12	TOKAQREQSURU	C-R71	
0239 12	TOMUKAQKEINI	C-R09	
0240 12	TOMUKAQKEINO	C-R71	
0241 12	TOONAZIKURAI	C-R52	
0242 12	TOUITUTEKINI	C-D10	
0243 12	TOHAKAGIRAZU	C-R69	
0244 12	TUGININOBERU	C-T00	
0245 12	UETODOUYOUNI	C-D10	
0246 12	UTINOIZUREKA	C-S48	K-71
0247 12	WOMOKUTEKINI	C-R08	
0248 12	WOMOKUTEKINO	C-R71	
0249 12	WOMOTONISITA	C-R71	
0250 12	WOMOTONISITE	C-R19	
0251 12	YOUNIOMOWARE	C-A44	K-71

図2 機械辞書の一部

#### 5. おまじ

本稿では、対象分野に余り依存しないと思われる言語情報を整理し、日常の日本語文の形態素解析を行う際に有効なルールの体系を示した。本研究のおもな特徴は、長単位の要素表現を広く取り入れて構造規定を行っており、意味処理との繋がりが直接的であるこ

と、およびかなり広範な言語表現に対処し得ることである。取扱う言語表現の範囲を広げることには、それだけ適格な意味の表現・伝達が可能となることであり、文解析でも、よりあいまいさの少ない意味解析の可能性をもたらす。無論、このためには、多様な表現について、使用される状況との関連も含めて、意味の整理を行っておく必要がある。我々は、即ち表現、関係表現の意味分類を行っているが、これについては別の機会に報告したい。

本稿で整理した言語情報は、資源としてまとめると、

- ・機械辞書(図2参照)
- ・接続ルール表(図1参照)
- ・活用語尾表(表3)
- ・音便変化文字表(表4)

(・活用接続条件表(表5))

であり、我々はすでにこれらを用いた形態素解析実験システムを作成し、要素表現の網目生、接続ルール(活用に関するルールを含む)の妥当性に関する検証を行って良好な結果を得ており、技術論文等の口語文章体文において体言、用言(および相当する表現)を除いた部分の構造は、ほとんど例外なく整理できたと考えている。

また、本稿の文法ルールは、要素表現の意味による、文中の隣接する場所での共起制限を含んでおり、接続を決定する局所レベルの解析で要素表現の多義の低減が可能を限り行える様に配慮されている点も確認済みである。例えば、「用言したため、…」における「ため」の意味は、先行する「た」との共起制限から目的の意味を取り得ず、原因の意味に限定できる。

接頭語、接尾語を含む複合名詞等の構造解析の問題、構文・意味解析への橋渡しの問題等、残された問題は多く、現在、検討中である。また、上記の資源に字種情報を加えて、ベテ書き文の

E-文節自動分析から書きに利用する実験を行っている。

## 参考文献

- (1)首藤、橋原、吉田：“日本語の機械処理のための文節構造モデル”，信学論(D)，J62-D，12 (1979)。
- (2)首藤公昭：“文節の構造と文解析”，電気学会連大(昭54)，(1979)。