

単語の釈義文を利用した単語間の階層関係の抽出について

鶴丸弘昭・水野浩司・内田 彰・日高 達・吉田 将
(長崎大学工学部) (九州大学工学部)

1. まえがき

近年、機械翻訳、自然言語理解など自然言語の機械処理に関する研究が盛んに行なわれるようになってきている。自然言語の機械処理において意味処理⁽⁵⁾を本格的に行なうためには、実用規模の意味辞書(広い意味でのソーラス)の開発が必要不可欠である。我々は、市販の国語辞典を高度に活用して、実用規模の意味辞書の開発に関する基礎的な研究を進めている^{(6),(7)}。ここでの基本的な問題の一つに、単語間の階層関係をどのようにして求めるかがある。

この問題に対する一つのアプローチとして、国語辞典の語釈義文(DS: Definition Sentence)の構造的特徴を利用して、見出し語(EW: Entry Word)に(階層)関係のある語(DW: Definition Word)を抽出し、DWとEWとの間の階層関係を抽出するシステム(階層関係付けシステム)の開発を試みている。語釈義文の構造的な特徴として、(a) DWは語釈義文の文末に表われる場合が多い⁽⁷⁾、(b) DSの文末にDWとEWとの関係を規定する表現(特殊文末表現)が含まれている場合がある、などがある。

本報告は、これらの特徴を利用した階層関係付けシステムの概要とその実験結果について述べたものである。

システムの構成を図1に示す。このシステムは、次の5つの処理からなっている。

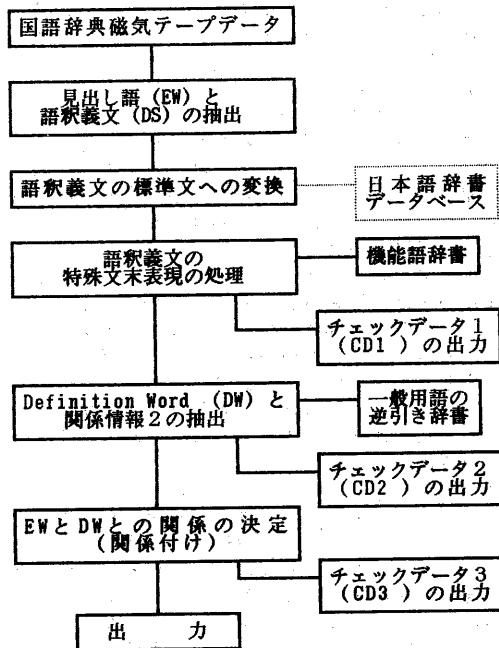


図1 システムの構成

(1) 見出し語 (EW) と語釈義文 (DS) の抽出。

国語辞典磁気テープデータから品詞別にEWに対応したDSを抽出する。一つのEWが複数個の語義を持っている場合があるので、それぞれの語義に対応したDSを抽出する必要がある。

(2) 語釈義文 (DS) の標準文への変換。

DSには、括弧やドット(・)による省略表現があったり、漢字に読みが付加されていたりするので、DSを標準的な日本語文表現に変換しておく。例【精鋭】: 勢いがよく鋭い・こと(兵士)。

標準文 { 1. 勢いよく鋭いこと。
 2. 勢いよく鋭い兵士。

(3) 語釈義文 (DS) の特殊文末表現の処理。

標準文変換されたDSから、Skeleton Sentence (SS)の抽出、およびEWとDWとの関係付けに必要な関係情報1(DSの特殊文末表現から得られる情報)の抽出を行なう。SSとは、DSが特殊文末表現を含まない場合はDSそれ自身であり、DSが特殊文末表現を含む場合はDSからそれを削除した残りの部分(文)である。

(4) Definition Word (DW) と関係情報 2 の抽出。

(3)で求められたSSから、DWの抽出、およびEWとDWとの関係付けに必要な関係情報2(SSの構造的な特徴から得られる情報)の抽出を行なう。

(5) EW と DW との関係の決定。

(3)と(4)で求められた関係情報をもとにEWとDWとの(階層)関係の決定(関係付け)を行なう。

先に、(1)と(2)を中心に発表した⁽⁴⁾ので、本稿では、(3)、(4)および(5)について報告する。

なお、本研究では、三省堂の新明解国語辞典磁気テープデータ⁽²⁾を使用している。

2. 語釈義文の特殊文末表現の種類と機能

標準文変換された語釈義文約10,000文を対象に、それらに含まれている特殊文末表現を調査して、約40種の(文末)機能パターン(FP: Functional Pattern)と約140種の(文末)機能語(FE: Functional Expression)を求めた⁽¹³⁾。FPは、特殊文末表現の構造を与えるもので、FEとの組合せで、DWとEWとの関係を規定する機能を持つ。ここでは、機能パターン(FP)を表1に示す型に整理している。

機能語(FE)の例を表2に示す。FEには、DWとEWとがどのような位相(見方)での(階層)関係かを明らかにするものがある。これらは、一般用語のソーラスにおいて有用な情報だと考えている。

特殊文末表現には、階層関係(上位-下位関係(>))、同義関係(=)、全体-部分関係(⊃)以外の関係を表わすものもある。このため、関連関係(R)を導入している。

表1. 語釈義文の機能パターンの例

型	機能パターン(FP)
100	「…DW」 + (など) + σ^e + FE.
200	…DW + (など) + の + FE.
201	…DW + (など) + W_{202} + の + FE.
202	…DW + (など) + と + DW + (など) + との + FE.
203	…DW + (など) + と + DW + (など) + との + FE.
300	…DW + (など) + を + σ^e + W_{301} + FE.
301	…DW + (など) + に対する + FE.
302	…DW + (など) + に対する + FE.
400	…DW + など.

注) σ^e : 任意の文字列

W_{202} : 意、の意、として、,、,、の中、の種、の一種

W_{301} : 言う、いう、言った、呼ぶ、呼んだ、指す、指した

表2. 機能語辞書の内容 (一部)

FE	k	s	FP	r1
一つ	2	2	201	>
			202	>
一つ一つ	2	1	201	》
群れ	2	1	201	R
横がわ	2	1	201	>
内がわ	2	1	201	>
表現	2	1	100	=>
雅語的表現	2	1	100	=>
漢語的表現	2	2	100	=>
			201	=>
称	1	3	100	=>
			201	=><
			202	=>
愛称	1	1	201	=>
異称	1	2	100	=>
			201	=>
言い方	3	2	301	=>
			100	=>
新しい言い方	3	1	100	=>

注) FE: 機能語, k: 文末語の文字数, s: 型の数
FP: 機能パターンの型, r1: DWとEWとの関係情報

3. 語釈義文の特殊文末表現の処理

3.1 機能語辞書

語釈義文 (DS) の文末表現が機能語 (FE) がどうかを判定するために、機能語辞書 (FED) を用いる。この判定は、DSの文末から最長一致法で行なわれる。このため、FEDは、拡張B-tree⁽¹⁰⁾の構造を持った逆引き辞書になっている。

FEDは、次の内容を持っている。

- (1) 逆向列の機能語 (FE)
- (2) FEの末尾の単語の文字数
- (3) FEと組合せ可能な機能パターン (FP) の個数
- (4) FEと組合せ可能なFPの型、およびそのFPが持っている関係情報1 (この情報は、DWとEWとの間にどのような (階層) 関係が存在するかを規定する)

表2に、FEDの内容の例を示す。

3.2 Skeleton Sentence と関係情報1の抽出

標準文変換された語釈義文 (DS) から、Skeleton Sentence (SS) と関係情報1の抽出を行なう。

処理手順は次の通りである。

Step 1. (機能語 (FE) の判定)

入力文 (標準文変換された語釈義文 (DS)、またはStep 4. からの再入力文) の文末表現がFEかどうかを、機能語辞書を用いて判定する。

- (1) 入力文の文末表現がFEでなければStep 5. へ。
- (2) 再入力された入力文の文末表現がFE (ただし、'など'を除く) であればStep 6. へ。

Step 2. (機能パターン (FP) の判定)

FP判定処理ルーチンにより、入力文のFPの型を決定する。FP判定処理ルーチンは、2. の表1をもとに作成している。入力文にFPが含まれていない場合は、Step 6. へ。

Step 3. (関係情報1の抽出)

機能語辞書よりFPに対応した関係情報1を抽出する。また、'など'が含まれていたことを関係情報1 ('など'情報) として抽出する。関係情報1が '?' であれば、Step 6. へ。

Step 4. (再入力文の生成)

入力文からFP (または 'など') を除去して、残りの部分を再入力文とする。Step 1. へ。

Step 5. (Skeleton Sentence (SS) の抽出)

Step 1. の(1)の条件を満足した文をSSとして抽出する。すなわち、入力文がFPを含む場合は、入力文よりFPを削除した残りの部分がSSとなる。入力文がFEを含まない場合は、入力文がそのままSSとなる。

Step 6. (チェックデータ1 (CD1) の出力)

入力文 (標準文変換されたDS) が、FEを含むがFPを含まない場合や、複数個のFEを含む場合などCD1として、DSがそのまま出力される。このような場合、現段階ではFEとDWとの区別が機械的には困難だからである。複数個のFEの組合せとしてどのような場合があるのか、そのときDWとEWとの関

係はどのようになるのか、などを調査するための資料としての利用を考えている。

図2に、SSおよび関係情報1の抽出例を示す。図2には、SSの一部、機能語辞書(FED)より求められた関係情報1、見出しとその正書法(EW)、語義の区別(大語義番号、小語義番号、標準文変換における通し番号)、および機能パターン(FP)の型が示されている。

図3に、CD1の出力例を示す。図3には、図2の項目の他に、チェックの理由が示されている。

4. Definition Word (DW)と関係情報2の抽出

3. で得られたSSから、Definition Word (DW)と関係情報2の抽出を行なう。

SSからDWを抽出するために一般用語(正書法の名詞見出し約75,000語)の逆引き辞書を用いた。この逆引き辞書は参考文献(12)で作成されたものである。

DWの抽出は、基本的には、この逆引き辞書を利用した、SSの文末との最長一致法によるマッチング処理である。しかし、語釈義文(DS)には特殊な記述形式(たとえば、動植物名は片仮名表記になっているなど)があるので、それらを考慮してDWおよび関係情報2の抽出を行なった。

処理手順は次の通りである。

Step 1. (漢字列または片仮名列の抽出)

入力文(SS、またはStep 6. からの再入力文)の文末に現われる漢字列または片仮名列を抽出する。抽出された文字列の文字数をnとする。

Step 2. (文末語の抽出)

一般用語の逆引き辞書とマッチングした、SSの文末語を抽出する。抽出された文末語の文字数をmとする。

Step 3. (チェックデータ2(CD2)の出力)

n = m = 0であればCD2としてSSが出力される。すなわち、SSの文末語が、片仮名列や漢字列でなく、逆引き辞書に登録されていない場合である。

Step 4. (DWの抽出)

a. n > mであればStep 1. で得られた文字列をDWとする。このとき、mをスタックしておく。
b. n ≤ mであればStep 2. で得られた文末語をDWとする。

Step 5. (関係情報2の抽出)

DWの直前に並列記号(‘、’、‘や’、‘と’)があるかどうかを判定する。

(1) 並列記号があれば、関係情報2として‘DWが複数個ある’情報を抽出する。と同時に、残りのDWを抽出するために、並列記号から文頭側の部分を再入力する(Step 1. へ)。

SS	r1	EW	a b c	FE	FPの型
バラ =>		いばら	(000)	雅語的表現	100
塩化ビニール =>		えんび【塩ビ】	(000)	略	100
盲人の官位 >		べつとう【别当】	(040)	第二	201
羽 R		うじよう【羽状】	(001)	形	201
形状・性質から見た、雲 R		うんきゆう【雲級】	(000)	分類	201
嫁 =>		よめご【嫁御】	(000)	尊敬語	302
昔の地方区画 >		どう【道】	(010)	一つ	201
もと、生物学・鉱物学・地質学 =><		はくぶつ【博物】	(000)	総称	201

注) a: 大語義番号, b: 小語義番号, c: 標準文変換処理における通し番号, r1: 関係情報1

図2. Skeleton Sentence (ss)と関係情報1の抽出例

DS	r1	EW	a b c	FE	FPの型	理由
文字を大きく二つに分けた一つ。		おんじ【音字】	(000)	一つ		1-1
寒の期間。		かんちゅう【寒中】	(000)	間		1-1
許された範囲の外。	?	らちがい【埒外】	(001)	外	201	2-1
下の部分。	?	かぶ【下部】	(010)	部分	201	2-1
接頭語と接尾語の総称。	=><	せつじ【接辞】	(000)	総称	201	2-2
生け花で、たけの低い草花類の総称。	=><	くさもの【草物】	(000)	総称	201	2-2

注) a, b, c, r1は、理由: 1-1 ...FPを満足しない
図2と同じ 2-1 ...FPを満足するが、FEの関係情報が‘?’
2-2 ...FPを満足するが、FEが複数個ある

図3. チェックデータ1 (CD1) の出力例

(2) 並列記号がなければ、SSが単語 (DW) かどうかを判定する。SSが単語であれば、関係情報2として 'SSが単語である' 情報を抽出する。

Step 6. (DWの出力)

Step 4. で抽出されたDWを出力する。

図4に、DWおよび関係情報2 (関係情報1も含む) の抽出例を示す。図3には、抽出されたDW (複数個の場合は '・' で区切られている) がKWIC形式で示されている。また関係情報1、見出しとその正書法 (EW)、語義の区別、機能パターン (FP) の型 (これらは図2と同じ)、および関係情報2が示されている。

図5に、CD2 の出力例を示す。

5. DWとEWとの関係の決定 (関係付け)

3. で求められた関係情報1、と4. で求められた関係情報2を用いて、DWとEWとの関係の決定 (関係付け) を行なう。

関係付けのための条件は次の通りである。

- i. DW ≧ EW (全体-部分関係)
関係情報1に '≧' 情報があり、それだけである。
- ii. DW R EW (関連関係)
関係情報1に 'R' 情報があり、それだけである。

iii. DW < EW (下位-上位関係)

- a. 関係情報1に 'など' 情報があり、かつDSがFPを含んでいる。
- b. 関係情報2に 'DWが複数個ある' 情報があり、かつDSがFPを含んでいる。

iv. DW = EW (同義関係)

- a. 関係情報1に '=' 情報があり、かつ関係情報2に 'SSが単語である' 情報がある。
- b. DSとSSとDWの3つが同じ場合、EWの定義が単語でのいいかえになっている。したがって、便宜的に、同義関係に近い関係として 'いいかえ関係 (≧)' を導入し、DW ≧ EW で表わす。

v. チェックデータ3 (CD3)

- a. 関係付けにあいまいさが残る。
- b. SSから複数個のDWが抽出されているが、DSがFPを含んでいない。(DW < EWにならない場合があるので)
- c. DWから文頭側の5文字以内に '・' がある。('複数個のDWがある' 情報の抽出ミスをチェックするため)

vi. DW > EW (上位-下位関係)

i ~ v の条件を満たさない。

図6に、EWとDWとの (階層) 関係付けの出力例を示す。図6には、抽出されたDW (KWIC形式で示されている)、決定された (階層) 関係、見出しとその正書法 (EW)、語義の区別、機能語 (FE)、機能パ

	DW	r1	EW	abc	FE	FPの型	d	e
	雄しべ・雌しべ	=><	ずい【蕊】	(000)	総称	201	4	
	気体・液体	<	えきたい【流体】	(000)	総称	203	4	
	股	》	うちもも【内股】	(000)	内側	201		*
盲人の	官位	>	べつとう【别当】	(040)	第二	201		
その人の犯した	犯罪	>	ざいめい【罪名】	(000)	名	201		
相手の	批評	=>	こうひ【高批】	(000)	敬称	201		
	終電車	=>	しゅうでん【終電】	(000)	略	100		*

注) a, b, c, r1は、図2と同じ
d: 'など' 情報として1, 'DWが複数である' 情報として4を割り当てたときの和
e: 'SSが単語である' 情報として*

図4. Definition Word (DW) と関係情報2の抽出例

	SS	r1	EW	abc	FE	FPの型
	相手によらず、人間愛をもつてつきあう	=>	ゆう【友】	(000)	漢語的表現	100
	小さなしろ・とりで	=>	じようるい【城堡】	(000)	漢語的表現	100
	しゆんば	=>	しゆんめ【駿馬】	(000)	古語的表現	100
	夜分になつてから	=>	ぼや【暮夜】	(002)	漢語的表現	100
	夕暮れに月が出ている	=>	ゆうづくよ【夕月夜】	(000)	雅語的表現	100
	いなか	=>	じかた【地方】	(010)	老人語	100
	他人・相手のむすこ	=>	れいそく【令息】	(000)	敬称	201
	ばらせん	=>	ばら	(020)	略	100

注) a, b, c, r1は、図2と同じ

図5. チェックデータ2 (CD2) の出力例

ターン (FP) の型、関係情報 1、関係情報 2 が示されている。

図 7 に、CD3 の出力例を示す。図 7 には、図 4 の項目の他に、チェックの理由が示されている。

6. 実験結果

本実験システムは、九州大学大型計算機センター FACOM M-382、長崎大学情報処理センター FACOM M-180 II AD 上に、PL/I および FORTRAN77 で実現されている。ここでは、3. の Skeleton Sentence (SS) と関係情報 1 の抽出、4. の Definition Word (DW) と関係情報 2 の抽出、および 5. の DW と EW との (階層) 関係の決定 (関係付け)、の実験結果について考察する。

実験用のデータとして、国語辞典磁気テープデータから 1/20 縮小辞典を作っている。1/20 縮小辞典には、単語が品詞別・活用別に国語辞典の 1/20 の割合で含まれている。単語の品詞・活用を決めるために、京都大学の調査資料^{(3),(4)}を利用した。この 1/20 縮小辞典から抽出した名詞見出し語の語釈義文のうち、標準文変換された語釈義文 (DS) 2,824 文を実験に用いた。

標準文変換された語釈義文 (DS: 2,824 文) に対して、SS が抽出された DS は 2,711 文 (約 96%) であり、チェックデータ (CD1) として出力された DS は 113 文 (約 4%) であった。

前段で抽出された SS (2,711 文) に対して、文末から単語が抽出された SS が 2,502 文 (約 92.3%)、チェックデータ 2 (CD2) として出力された SS が 209 文 (約 7.7%) であった。前者 (2,502 文) の

	DW	関係	EW	a b c	FE	FPの型	r1	d	e
ているものから出る	顔	》	よこつつら【横っ面】	(010)	横がわ	201	》		*
	一年	》	はんとし【半年】	(001)	半分	201	》		
	木	R	じゆか【樹下】	(000)	下	201	R		*
	炎・熱・光	<	ひ【火】	(012)	など	400	<	5	
	気体・液体	<	りゆうたい【流体】	(000)	総称	203	<	4	
	イノシシ	=	やちよ【野猪】	(000)	漢語的表現	201	=>		*
	追試験	=	ついし【追試】	(020)	略	100	=>		*
	植物の	ストック	>	あらせいとう	(000)	和名	201	=>	
	相手の	元素	>	しゅうそ【臭素】	(000)	一つ	201	>	*
		説	>	こうせつ【高説】	(000)	敬称	201	=>	
運送の 秋の澄み渡つた 物に乗らず、足で歩く 扱い方が困難な	イギリス	≧	えいこく【英国】	(000)					*
	和歌	≧	うた【歌】	(020)					*
	名前	≧	めいしやう【名称】	(000)					*
	料金	>	うんちん【運賃】	(000)					
	空	>	あきぞら【秋空】	(000)					
	こと	>	とほ【徒歩】	(000)					
	もの	>	なんぶつ【雜物】	(001)					

注) a, b, c, d, e, r1 は、図 4 と同じ

図 6. 階層関係システムの出力結果の例

	DW	r1	EW	a b c	FE	FPの型	d	g
目的として飼う、	鶏	><	らんようしゆ【卵用種】	(001)	品種	201		0
平安時代の	皇后	=><	ちゆうぐう【中宮】	(001)	称	201		0
それぞれの	省	》R	しやうない【省内】	(001)	内部	201		0
題詠でなく作った	詩・歌		むだい【無題】	(010)			4	2
農業用の	機械・器具		のうきぐ【農機具】	(000)			4	2
	正号・負号		せいふ【正負】	(010)			4	2
隣の隣国・あたり	近所	=>	しりん【四隣】	(001)	漢語的表現	100		3
に記入・計算する	ところ		こうざ【口座】	(020)				3
深い川・海へ移る	魚		おちうお【落(ち)魚】	(020)				3

注) a, b, c, d, r1 は、
図 4 と同じ

g はチェックの理由: 0...関係決定にあいまいさが残る
2...FPがなく、DWが複数ある
3...DWから文頭側 5 文字以内に '・' がある

図 7. チェックデータ 3 (CD3) の出力例

中で、DWが正しく抽出されていた SS が 2,386文 (約95.4%)であった。

前段でDWが正しく抽出されたSS (2,386 文) に対して、関係情報1と関係情報2を用いて関係付けがなされたDW-EWの組が2,305 個 (96.6%)、CD3が81個 (3.4%)であった。前者 (2,305 個) の中で、妥当な関係付けは2,297 個 (99.6%)であった。

以下、各段階での実験結果について考察する。

(I) Skeleton Sentence (SS)、関係情報1、
チェックデータ1 (CD1)

SSが抽出された語釈義文 (DS) 2,711 文に対して、(1)機能語 (FE) を含まないDSが2,374 文 (約87.6%) (2)FPを含んでいるDSが337 文 (約12.4%)であった。

(2)の、337 文の内訳は次の通りである。

100型	214文	201型	113文
200型	114文	203型	1文
300型	3文	301型	2文
400型	6文	302型	1文

CD1 (113 文) の内訳は次の通りである。

a. FEがありFPを満足しないもの	100 文
b. FEがありFPを満足するもので、	
b-1. FEの関係情報が '?' である	10文
b-2. FEが複数個ある	3文

(II) Definition Word (DW)、関係情報2、
チェックデータ2 (CD2)

文末語が抽出されたSS (2,502 文) に対して、(1)SSの文末語として正しい語が抽出されているSSが2,406 文 (96.2%)、(2)文末語として誤った語が抽出されたSSが96文 (3.8%)であった。

(1)、(2)の内訳は次の通りである。

(1)文末語として正しい語が抽出	2,406個
i. 文末語がDWとして妥当	2,386個
ii. 文末語がDWとして不適当	20個
a. 文の記述または構造による	8 個
a-1. DWとして不都合	6 個
a-2. 抽出したFEがDW	2 個
b. 標準文変換ができていない	12個
b-1. ルビが取れていない	5 個
b-2. 正書法が () の中にある	2 個
c. 1 個のDWしか抽出していない	5 個
(2)文末語として誤った語が抽出	96個
i. 品詞の誤り	17個
a. 見出しが名詞でない	7 個
a-1. 処理過程で品詞決定を誤った	2 個
a-2. 国語辞典中でおかしい	5 個
b. 見出しは名詞であるが文末語が名詞でない	4 個
c. 子見出しである	6 個

ii. 逆引き辞書の見出しの不備	72個
a. 文末語が正書法で書かれているが、逆引き辞書に登録されていない	13個
a-1. 平仮名書き	3 個
a-2. 平仮名・漢字混合	9 個
a-3. 片仮名・漢字混合	1 個
b. 文末語が正書法で書かれていない	56個
b-1. 平仮名書き	44個
b-2. 平仮名・漢字混合	12個
c. 逆引き辞書との最長一致によるミスマッチング	3 個
iii. 標準文変換ができていない	3 個
a. ルビが取れていない	1 個
b. 正書法が () の中にある	1 個
c. その他	1 個
iv. データエラー	4 個

(1)、(2)の内訳の例を次に示す。a-1.がa1. に対応している。表記は、図4の一部が省略された形式になっている。ただし、関係情報1をDWとEWとの間に入れていない。(1)のi. の例は、図4に示されている。

(1)のii. の例

a1. …あたる 年・月・日【甲子】(010)	
a2. 備えるべき、一定【体裁】(031) 形	201
b1. 主人と客 キャク【主客】(011)	
b2. …ぐあい 具合【ぐわい】(002)	
c. 高い所と低い 所【高低】(001)	

(2)のi. の例

a1. …ことを表わす【逆】(200)	
a2. …ように美しい【花の】(010)	
b. …そうでないか【真否】(000)	
c. …ことを表わす【何かと言えば】(000)	

(2)のii. の例

a1. …おもちゃ【おしやぶり】(000)	
a2. …浮彫り【彫上げ】(000)	
a3. …塩化ビニール【塩ビ】(000) 略	100
b1. …表わすくらゐ【段位】(000)	
b2. …人の顔つき【人相】(001)	
c. …として払う金【場銭】(020)	

(2)のiii. の例

a. 行在アンザイ 所【行宮】(000)	
b. …かし・仮借【仮借】(022)	
c. そば 屋一品料理【2】【店屋物】(002)	
d. …芸術の 行動【芸術祭】(000)	

CD2 として出力されたSSは、SSの文末単語が片仮名列や漢字列でなく、逆引き辞書に登録されていないものである。

国語辞典⁽¹⁾の語釈義文に用いられている単語は必ずしも漢字表記の正書法でない (平仮名書きなど) 場合があるので、機械処理を困難にしている。今回、利用した逆引き辞書が正書法のみから構成さ

れていることも、DWが誤って抽出されたり、SSがCD2として出力されたりする原因の一つになっていると考えられる。

(2)のi.のa-1.は、本システムの第1段階(見出し語とその語釈義文の抽出)での誤りである。また、(2)のi.のa-2.は、国語辞典⁽¹⁾の記述形式のあいまいさによるもので、機械的に誤りをなくすことは困難である。

(1)のii.のb.は、本システムの第2段階(標準文変換)での誤りであり、関係情報2の抽出失敗の原因にもなっている。

(III) DW-EW関係付け、チェックデータ3 (CD3)

DWとEWとの関係付けがなされたDW-EWの組(2,317個)の内訳は次の通りである。

- | | |
|---------------------|---------|
| (1) 関係付けが妥当 | 2,297 個 |
| (2) 関係付けが妥当かどうか判定困難 | 6 個 |
| (3) 関係付けが誤り | 2 個 |

(2)、(3)のデータを次に示す。(1)の例は、図6に示されている。

(2)のデータ

歌舞伎・小説 < 【時代物】(021) など 400 > <
 悲しみ・喜び < 【哀歓】(001) 100 = >
 一步の長さ > 【半歩】(002) 半分 201 >
 離れた都市 > 【いなか】(012)
 き叫ぶ習俗 > 【哀号】(001)
 元旦の朝 > 【元旦】(002)

(3)のデータ

耳 R 【耳たぶ】(001) 部分 100 R
 つづり 字 > 【字母】(011) -つ- 201 >

(2)の、関係付けが妥当かどうか直感的に判定困難な場合でも、多義情報を考慮すれば妥当とみなせるものがある。(3)の、関係付けの誤りは機能語辞書の不備によるものであった。

CD3 (81個)の内訳は次の通りである。

- | | |
|--------------------------------|-----|
| a. 関係決定にあいまいさが残る。 | 17個 |
| b. 複数個のDWが抽出されているが、DSがFPを含まない。 | 49個 |
| c. DWから文頭側の5文字以内に'.'がある。 | 15個 |

b.の条件は、EW > DWの決定に利用しようとしたものもあったが、例外や判定困難な例がいくつか見出された。また、c.の条件は、'複数個のDWがある'情報の抽出ミスをチェックするためである。c.の条件で5文字で良いのかどうかなど、現在、CD3の出力データを検討している。

7. あとがき

本稿では、国語辞典の語釈義文(DS)の構造的特徴を十分に活用して、単語間の意味的關係の一つである階層関係を機械的に求める階層関係付けシステムの概要およびその実験結果について報告した。

実験結果の考察により、次のことが言える。

- (1)本システムは人間支援型を前提にしているが、入力文(標準文変換されたDS)2,824文に対して、チェックデータの総数は403文(14.3%)であった。抽出プログラムの改良のため、チェックデータを検討している。
- (2)DWの抽出に関して、文末語が抽出されたSS(2,502文)に対して、適切なDWが抽出されたSS(2,386文)の割合は95.4%であるが、正しい文末語が抽出されたSS(2,406文)に対しては、その割合が99.2%になる。
- (3)DW-EWの関係付けにおいて、適切なDWに対して、正しく関係付けられた割合は99.6%である。DWが正しく抽出されれば、関係付けはうまくいっているとみなせる。直感的に妥当かどうか判定困難な関係付けも、多義を考慮すれば、妥当とみなせるものがある。(III)の(2)のデータ参照)
- (4)DWとEWとの関係に、いかに関係(≧)を導入した。調査の範囲ではほとんどDW=EWとみなしてよかったが、DW>EWの可能性を考慮したためである。
- (5)一般に、部分-全体関係は、a.地理的位置関係、b.体の器官や組織、c.知識学問の領域の3つの場合に限定されている⁽¹⁵⁾が、時間的關係、距離的關係も、部分-全体関係とみなしてよいと思われる例があった。(図6参照)
- (6)DWとして'もの'や'こと'が抽出されている場合が多い。DWが正しく抽出されたSS(2,386文)に対して、文末が'こと'となるSSが422文(17.7%)、文末が'もの'となるSSが84文(3.5%)であった。
- (7)DWの抽出およびDW-EWの関係付けの精度をあげるためには、(a)逆引き辞書の整備・補強、(b)DSの構造的特徴のより細かい調査が必要である。

今後の研究課題として、次のようなものがある。

- (1)国語辞典に含まれる全ての名詞見出し語についてDWとEWとの(階層)関係付けを行なう。
- (2)階層関係にあるDWとEWとを利用して、階層構造を求める。
これについては、別の機会に報告する予定である。これは、階層構造が国語辞典の中に暗黙の形で埋め込まれているとの仮定のもとに、この階層構造を、陽な形で抽出しようとするものであり、シソーラスの自動作成の手がかりを与えるものであると考えている。この場合、単語の多義というやっかいな問題が係わってくる。
- (3)語釈義文に含まれている関係語(表現)⁽¹⁶⁾の抽出と調査。
- (4)語釈義文の、格情報を用いた係り受け解析⁽¹¹⁾を行ない、意味辞書作成支援システムの開発を試みる。

本稿では、名詞について考察したが、意味辞書では、動詞、形容詞、副詞等の意味の定義⁽⁸⁾も必要である。

謝辞

本研究の当初より種々の助言をいただいた九州大学工学部田中武美助手、福岡大学工学部吉村賢治講師に感謝致します。

参考文献

- (1) 金田一京助、金田一春彦、見坊豪紀、柴田武、山田忠雄(編)：新明解国語辞典、三省堂、第2版(1974)、第3版(1981)
- (2) 横山晶一：国語辞典データベース化の準備、電子技術総合研究所彙報、Vol.41、No.11、PP.19-27(1977.11)
- (3) 長尾真、辻井潤一、山上明、建部周二：国語辞典の記憶と日本語文の自動分割、情報処理、Vol.19、No.6、PP.514-521(1978.6)
- (4) 長尾真：計算機による日本語文章の解析に関する研究、昭和53年度文部省科学研究費特定研究(1)昭和53年度研究報告書(1979.2)
- (5) 田中穂積：計算機による自然言語の意味処理に関する研究、電子技術総合研究所研究報告、第797号(1979.7)
- (6) 首藤公昭：文節構造モデルによる日本語の機械処理に関する研究、福岡大学研究所報、第45号(1980.3)
- (7) 中野洋：分類番号つけ支援システム、情報処理学会計算言語学研究会資料25-5(1981.2)
- (8) S.Yoshida：CONCEPTUAL TAXONOMY FOR NATURAL LANGUAGE PROCESSING, in Kitagawa, I.(ed), Computer Science & Technologies, Ohmsha, Japan (1982)
- (9) S.Yoshida, H.Tsurumaru, I.Hitaka：MAN-ASSISTED MACHINE CONSTRUCTION OF A SEMANTIC DICTIONARY FOR NATURAL LANGUAGE PROCESSING, Proc. COLING82, PP.419-424(1982.7)
- (10) 日高達、吉田将、稻永紘之：拡張B-treeによる日本語単語辞書の作成、情報処理学会自然言語処理研究会資料33-8(1982.10)
- (11) 日高達、吉田将：格文法による日本語の構文解析、情報処理学会自然言語処理技術シンポジウム論文集、PP.41-PP.46(1983.6)
- (12) 吉村賢治、山下明男、日高達、吉田将：専門用語の自動収集システムについて、情報処理学会自然言語処理研究会資料42-1(1984.3)
- (13) 鶴丸弘昭、日高達、吉田将：国語辞典を利用したシソーラスの作成について、情報処理学会第27回(昭和58年度後期)全国大会、2H-2(1983.10)
- (14) 鶴丸弘昭、内田彰、日高達、吉田将：国語辞典からの情報抽出とその構造化、情報処理学会自然言語処理研究会資料43-6(1984.5)
- (15) 長尾真監修：日本語情報処理、電子通信学会(1984.5)