

自然語意味情報を対象にした 格構造データベース

*柴崎 真人 *小口 琢夫 **木内 伊都子

* (株)日立製作所 システム開発研究所

**(株)日立製作所 中央研究所

文書検索システム等の情報検索システムにおいて、蓄積される情報の内容や検索要求の表現にはキーワードのみでは不十分であり、より高度な情報表現が必要である。そこで、文書内容や関連知識などの自然語意味情報を、多ソート論理式で意味付けした格構造モデルで保持する演繹データベースを開発中である。本システムは、対象とする知識を自然語意味情報に限定し、そのための表現機能として、多ソート論理式のルール節の拡張にあたる包含連体修飾というデータ形式と、ソートの集合演算にあたる合成ソートを提供している。

A Case Structured Database for Natural Language Information

*Masato Shibasaki, *Takuo Koguchi, **Itsuko Kiuchi

*Systems Development Laboratory, Hitachi, Ltd.

**Central Research Laboratory, Hitachi, Ltd.

*1099 Ohzenji Asao-ku Kawasaki-shi, 215 JAPAN

**1-280 Higashikoigakubo Kokubunji-shi, 185 JAPAN

In order to represent the contents of documents or retrieval requests for information retrieval systems such as document retrieval systems, it has been recognized that conventional keyword systems are insufficient and a more sophisticated method is required. A deductive database system, which is based on many sorted logic and case structures of natural language sentences representing the contents of documents and retrieval requests, is being developed. This system features its representation capability of extended rule clauses and compound sorts.

1. はじめに

現在、文書検索システムの文書内容に関する検索方式としては、商用のものも、研究レベルでも、キーワード方式が中心である[1]。それに対し、キーワードのみでは、蓄積された文書の内容も、ユーザの検索要求も十分に表現できないということで、文書の内容と質問文を意味構造モデルで表現し、両者のマッチングをとることにより、質問文に適した文書を検索する方式の研究もすすめられている[2][3][4]。しかし、これらは、明確に意味付けされた意味構造モデルをもつものではない。我々は、文書内容・関連知識などの自然語意味情報を、多ソート論理式で意味付けした格構造モデルとして保持し、ロジックベースの健全で完全な検索処理を行う「知的検索システム」を現在開発中である。本システムにより、次のようなメリットが得られるものと考えている。

- (1)文書内容や書誌事項に関する知識は、検索対象とするだけでなく、検索時に利用することもできる。
- (2)検索文内の概念をより詳細な情報に置き換えることによって、検索結果を提供できるので、リファレンスDBとしてだけでなく、ファクトDBとしての機能も果たす。

本稿では、本システムのデータベース部分を格構造データベース(Case Structured Data Base)とよび、このデータベースモデルを中心に報告する。ビジュアルインターフェース部分については[5]を参照のこと。

2. 基本的考え方

インスタンス(個別の事物)のみでなく、クラス(個別の事物を包含する一般的な概念)も記述の対象として含むことにより、継承階層の表現・検索を効率的に実現しようとする拡張Prologがいくつか提案されている(以後、継承階層Prologとよぶ)[6][7][8][9]。格構造データベースのデータベースモデルおよび検索処理の单一化手順はこれらの考え方を多分に継承しているが、次の点が大きく異なる。

- (1)各述語の各々の引数が自然語のいずれの格に相当するかを引数の位置によってではなく、格と引数の対として記述することにより示す。
 - (2)述語にも階層構造をもたせ、また、述語の引数として時制という要素ももたせる。
 - (3)接合連体・包含連体修飾というデータ形式を提供する。
 - (4)格の中身は単一のソートまたは定数ばかりではなく、それを集合演算した合成ソートの記述も許す。
 - (5)変数同志の单一化の際に、変数を型付けしている合成ソートが互いに包含関係にはないが、共通部分集合をもつとき、单一化代入として、変数を型付けしている合成ソート間のglb(greatest lower bound)を直接求めるのではなく、単に、合成ソートを "▷" でつなげる。
 - (6)ソートと定数の階層構造は、包摂(is_a)関係のみではなく、部分全体(part_of)関係も含み、検索の際には、部分全体関係も利用する。
 - (7)単位を伴う単一の数値または1組の数値も定数・ソートと同様に扱う。
 - (8)定数・ソート・述語は内部的にはIDで管理する。これにより、同義語のカテゴリ化や同義異義語の区別が可能である。
- (1)・(2)は本章で、(3)は次章で、(4)・(5)は4章で、(6)・(7)・(8)については5章で説明する。
なお、本稿では、分りやすさを主眼とすることにし、正確な定義は省略し、具体例を中心に説明する。

2. 1 概念階層知識

格構造データベースが扱う知識には、概念階層知識と文知識とテンプレートがある。

概念階層知識は、定数・ソートと述語の階層構造を示す。図1に例を示す。ソートには "@" をつけています。我々は、現実的な観点から、多重継承を許すことにして、また、階層構造が束(lattice)であることを仮定しない。以降の例では図1の概念階層を仮定する。

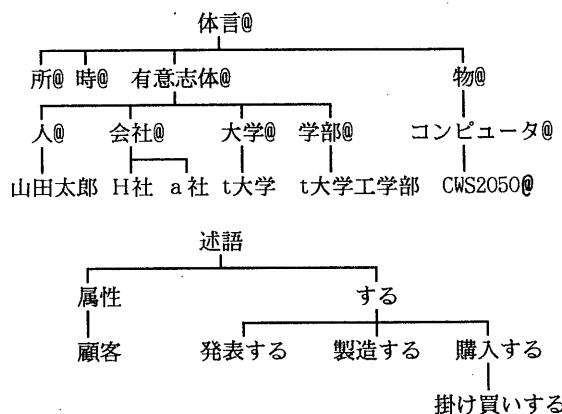


図1 概念階層の例

表1. C S D B の格ラベル一覧

略記コード	英語名	日本語名
A	Agent	動作主(何が)
O	Object	対象(何を)
R	Recipient	受益者(何に)
S	Start	起点(何処から)
D	Destination	方向(何処へ)
P	Place	所(何処で)
T	Time	時(いつ)
C	Cause	原因(なぜ)
I	Instrument	手段(何で)
U	Purpose	目的(何のために)
M	Measure	数量(いくつ/どの位)
E	Entity	属性対象(～の…は)
V	Value	属性値(である)

2. 2 文知識

文知識は、単文・言い換え規則・接合連体修飾・包含連体修飾からなる。本節では、単文と言い換え規則について説明する。接合連体修飾・包含連体修飾については次章で述べる。単文は継承階層Prologのファクト節に、言い換え規則はルール節に対応する。

「a社はあるコンピュータをすべての在京新聞で発表した」という内容を考える。これは、多ソート論理式では、

$$\exists Xc / \text{コンピュータ}, \forall Xn / \text{在京新聞} \quad \text{発表した} (a\text{社}, Xc, Xn)$$

と表せる。このとき、コンピュータと在京新聞はソートである。この式をスコーレム化すると、

$$\forall Xn / \text{在京新聞} \quad \text{発表した} (a\text{社}, \text{コンピュータ} \# 1, Xn)$$

となる。このように意味づけされたものを格構造データベースでは、

$$<\text{発表する, 過去, } a\text{社} / A, \text{コンピュータ} \# 1 / O, \text{在京新聞} @ / I>$$

という形式で表現する。このタイプの知識を単文という。コンピュータ # 1 は、システムが内部的につくりだしたコンピュータのソートに属するスコーレム定数であり、ビジュアルインターフェースを介する

ことにより、ユーザからは「あるコンピュータ」としか見えない。一番左が述語、その次が時制である。時制は、現在・過去・未来のうちのいずれかである。そのあとは、格とその中身との対応を示している。格構造データベースでは、表1に示す13種類の格を用意しており、このうちの使用する格に対してのみ、対応する中身を記述する。格の出現順序に意味はない。

本来、

$\forall X_n / \text{在京新聞}, \exists X_c / \text{コンピュータ} \text{ 発表した } (a \text{ 社}, X_c, X_n)$

(「a社はすべての在京新聞で何らかのコンピュータを発表した」)

という多ソート論理式も表現できるべきだが、(1)ユーザに限量子の依存関係を指摘させるのは困難、(2)スコーレム関数の導入に伴う出現試験(occur check)の計算の手間を省くという理由から現バージョンでは、前置する限量子が $\exists \dots \exists \forall \dots \forall$ という順序の多ソート論理式のみを対象にしている。

次に、「ある会社が製造した物を購入した人を、その会社の顧客という」という内容を考える。これは多ソート論理式では、

$\forall X_c / \text{会社}, \forall X_h / \text{人}, \forall X_t / \text{物}$

製造した(X_c, X_t) \wedge 購入した(X_h, X_t) \rightarrow 顧客(X_c, X_h)

と表せる。このように意味づけされたものを格構造データベースでは、

<顧客、現在、会社@1/E, 人@2/V>

← <製造する、過去、会社@1/A, 物@3/O>

<購入する、過去、人@2/A, 物@3/O>

S 1

という形式で表現する。このタイプの知識を言い換え規則という。会社@1の@1は、2つの会社@1が同一事物であることを意味するものであり、システムが自動的に番号付けして付与する。

2. 3 テンプレート

テンプレートは各々の述語に対して、とりえる格とその型付けを示す。これは下位の述語にも継承される。例えば、購入するという述語のテンプレートは、

<購入する、有意志体/A, 物/O, 所/P, 時/T, 物/U, >

のように表現されている。システムは、テンプレートに適合しない文知識や質問文はうけつけない。

2. 4 質問文

質問文は継承階層Prologのゴール節に相当する。「H社の顧客である人は?」という質問を考える。「H社の顧客は人である」という内容は、多ソート論理式では、

$\exists Y_h / \text{人} \text{ 顧客}(H \text{ 社}, Y_h)$

と表わされ、これの否定は、

$\forall Y_h / \text{人} \sim \text{顧客}(H \text{ 社}, Y_h)$

である。格構造データベースでは、先の質問文を、

← <顧客、現在、H社/E, 人@/V>

という形式で表現する。

2. 5 単一化

以後は Prolog の用語に準じて、单一化方法を説明する。選択式と入力節とが单一化に成功する基準は以下のすべての条件を満足するときのみである。

(1) 選択式の述語が入力節のヘッド部の述語と同じか、より上位である。

(2) 選択式の時制と入力節のヘッド部の時制とが等しい。

(3) 選択式に記されているすべての格に対して、次の条件を満足する。

選択式の格と同一の格が入力節のヘッド部に存在し、ソートをそのソートに属する定数の集合、定数を単位集合とみなしたとき、両者の格に対応する定数またはソートが包摂関係にあるか、共通部分集合をもつ。

言い換え規則 S 1 と、S 2・S 3 の單文があるとき、「H社の顧客は？」という質問に対する検索処理の一例を示す。

<製造する, 過去, H社/A, CWS2050@/O>

S 2

(「H社はすべてのCWS2050を製造した」)

<掛け買いする, 過去, 山田太郎/A, CWS2050#/1/O>

S 3

(「山田太郎はあるCWS2050を掛け買いした」)

最初のゴール節は、

← <顧客, 現在, H社/E, 人@/V>

G 1

である。このゴール節はまず、S 1 を入力節として、代入を {会社@1 → H社/E, 人@ → 人@2/V} とすることにより、单一化に成功し、新たなゴール節

← <製造する, 過去, H社/A, 物@3/O>

G 2

<購入する, 過去, 人@2/A, 物@3/O>

G 3

を得る。次に、G 2 を選択式とするとき、S 2 を入力節として、代入を {物@3 → CWS2050@/O} とすることにより、单一化に成功し、新たなゴール節

← <購入する, 過去, 人@2/A, CWS2050@/O>

G 4

を得る。最後に、選択式 G 4 に対して、S 1 を入力節として、代入を {人@2 → 山田太郎/A} とすることにより、单一化に成功し、空節が導かれ、G 1 の人@/Vに対する合成代入 {人@ → 山田太郎/V} が解となる。

3. 連体修飾

我々は、自然語意味情報をより直接的に表現するために、接合連体修飾と包含連体修飾という 2 種類のデータ形式を提供する。

3. 1 接合連体修飾

「a 社が製造したコンピュータを t 大学が購入した」という内容を考える。これは、「a 社が製造した」コンピュータの集合と、「t 大学が購入した」コンピュータの集合との間に共通部分があることを意味する。多ソート論理式では、

$\exists Xc/\text{コンピュータ 製造した}(a\text{ 社}, Xc) \wedge \text{購入した}(t\text{ 大学}, Xc)$
と表わされ、スコーレム化すると、

製造した(a 社, コンピュータ # 1) \wedge 購入した(t 大学, コンピュータ # 1)
となる。格構造データベースでは、2つの單文
<製造する, 過去, a 社/A, コンピュータ # 1/O>
<購入する, 過去, t 大学/A, コンピュータ # 1/O>
の合成として表現する。

3. 2 包含連体修飾

「t 大学が購入したコンピュータはすべて a 社が製造した」という内容を考える。これは、「a 社が製造した」コンピュータの集合の中に、「t 大学が購入した」コンピュータの集合が含まれることを意味する。「t 大学が購入した」コンピュータが空ではないと考えるなら、多ソート論理式では、

$$(\forall Xc/\text{コンピュータ 購入した}(t\text{ 大学}, Xc) \rightarrow \text{製造した}(a\text{ 社}, Xc)) \wedge \\ (\exists Xc/\text{コンピュータ 購入した}(t\text{ 大学}, Xc))$$

と表わされる。これをそのまま格構造データベースで表現するなら、

$$<\text{製造する, 過去, } a\text{ 社}/A, \text{ コンピュータ}@1/O> \\ \leftarrow <\text{購入する, 過去, } t\text{ 大学}/A, \text{ コンピュータ}@1/O> \quad S 4$$

$$<\text{購入する, 過去, } t\text{ 大学}/A, \text{ コンピュータ}#1/O> \quad S 5$$

となるが、格構造データベースでは、この他に、S 4 と S 5 から導かれる S 6 もデータとして保持しておく。

$$<\text{製造する, 過去, } a\text{ 社}/A, \text{ コンピュータ}#1/O> \quad S 6$$

ただし、S 4 が入力節として選ばれて单一化に成功したならば、S 5 と S 6 が入力節として選ばれることがないようにこれらにマスクをかけ、入力節としての S 4 の処理が終われば、マスクをはずす。

これにより、演繹推論を行わずに、「a 社がコンピュータを製造した」ことを導くことができる。

4. 合成ソート

格構造データベースでは、格の中身は单一の定数またはソートのみではなく、次のような構造(合成ソートとよぶ)を許している。ここで、g_{ij}はソート、c_{ij}は定数またはソートである。(i, j, n₁, ..., n_mは整数)

$$(g_{11} - (c_{12} + \dots + c_{1n_1})) \cap \dots \cap (g_{m1} - (c_{m2} + \dots + c_{mn_m}))$$

合成ソートは、ソートをそのソートに属する定数の集合、定数を単位集合とみなして、それらに集合演算を施したものである。

单一化のルールは次のとおり。

格 Cにおいて、合成ソート e₁, e₂に対して、

① e₁のすべての要素が e₂に含まれるとき

单一化に成功し、单一化代入は {e₂ → e₁/C} となる。

② e_2 のすべての要素が e_1 に含まれるとき

单一化に成功し、单一化代入は $\{e_1 \rightarrow e_2 / C\}$ となる。

③ 上のどちらでもなく、かつ、共通部分集合をもつならば、

单一化に成功し、单一化代入は $\{e_1 \rightarrow e_1 \cap e_2 / C, e_2 \rightarrow e_1 \cap e_2 / C\}$ となる。

④ 上のいずれでもないならば、单一化に失敗する。

5. 部分全体関係・数値情報・同義語

5. 1 部分全体関係

格構造データベースでは、定数とソートの階層構造を構成する関係として、包摂関係だけでなく、部分全体関係も利用している。これは、部分全体関係を利用した次のような検索が有用であると考えるからである。

<購入する, 過去, t 大学工学部 / A, CWS2050 # 1 / O >

S 7

(t 大学工学部が CWS2050 を購入した)

S 7 の單文と、 t 大学工学部が t 大学の部分であるという部分全体関係があるとき、「CWS2050を購入した大学は?」という質問に対して「 t 大学」が解として期待される。

検索処理の過程では、上記の質問を表すゴール節、

← <購入する, 過去, 大学 @ / A, CWS2050 @ / O >

を選択式とし、S 7 を入力節とする。このとき、代入 (大学 @ → t 大学 / A, t 大学工学部 → t 大学 / A, CWS2050 @ → CWS2050 # 1) により、单一化に成功すればよい。2.5 節で述べた、单一化の成功の基準(3)に次をつけ加える。

選択式の格と同一の格が入力節のヘッド部に存在し、そのヘッド部の格の中身が定数であり、その部分全体関係による上位概念（全体概念）で、選択式の格の合成ソートに包摂される定数 w が存在する。このとき、格を C 、入力節のヘッド部の格の定数を t 、選択式の格の合成ソートを e とするとき、单一化代入は $\{t \rightarrow w / C, e \rightarrow w / C\}$ 。ただし、入力節のヘッド部の述語と格の組合せが、本項目の適用を受けるものと定義されている場合に限る。

5. 2 数値情報

格構造データベースでは、文知識あるいは質問文において格に数値情報を表現することができる。数値情報は、単一の数値または数値範囲と、数値の単位からなる。数値情報における限量子の意味は、 \forall については、その数値範囲に属するすべての値についてその文知識（あるいは質問文）の内容が成り立つことであり、 \exists については、その数値範囲に属するある値についてその文知識（あるいは質問文）の内容が成り立つことである。次に例を示す。

<購入する, 過去, t 大学工学部 / A, [1985,1990]年 # 2 / T, CWS2050 # 1 / O >

(t 大学工学部は 1985年から1990年の間に CWS2050 を購入した)

5. 3 同義語

格構造データベースでは、定数・ソート・述語に対して概念IDを付与し、概念階層知識・文知識・テンプレート・検索文においては、概念IDによりそれぞれを識別している。それぞれの概念IDとそれに対応する文字列との対応は別に管理する。同義語には同一の概念IDを与えておくことにより、表記が異なっていても同一視することができる。また、同級異義語については、それぞれの意味ごとに概念IDを与えることにより、区別して扱うことができる。同様の考えは、[10]にもみられる。

6. おわりに

格構造データベースは、自然語意味情報を表現しやすくするために、連体修飾・合成ソート・部分全体関係・数値情報・IDによる概念の管理を導入した格構造データベースモデルをもつ演繹データベースである。本データベースは、リファレンスDBとしてだけでなく、ファクトDBとしての機能も果たすものであり、より高度な検索機能をもつ情報検索システムの要素技術になりうるものと考えている。

- [1] 拜原正人：日本語文献データベースへの知的アクセス：電子情報 信学会誌，72,7, pp. 797-806, (1989).
- [2] 杉山健司, 秋山幸司, 伊吹潤, 川崎正博, 内田裕士：自然言語理解に基づく情報検索システム IRIS：情報処理学会, 自然言語処理研究会資料, 58-8,(1986).
- [3] 稲垣博人：文書内容検索システムにおける現状と課題 一文書入力型内容検索システムを中心にして：情報管理, 33,2, pp. 123-135,(1990).
- [4] 福永博信, 齋藤珠喜：全文探索と多様な表現：情報処理学会第39回全国大会, 1G-7, pp. 682-683,(1989).
- [5] Fujisawa, H., Kiuchi, I., Koguchi, T., Kondo, H.: A Visual Interface for a Personal Information Base using a Concept Network: Proc. of The second International Symposium on Database Systems for Advanced Applications, Tokyo, (1991).
- [6] 赤間清：PAL：継承階層を扱う拡張Prolog：情報処理学会論文誌, 28,4, pp. 322-329,(1987).
- [7] Huber, M., Varsek, I.: Extended Prolog for Order-Sorted Resolution: Proc. of '87 International Symposium on Logic Programming, pp. 34-43,(1987).
- [8] Montini, G.: Efficiency Considerations on Built-in Taxonomic Reasoning in Prolog: Proc. of IJCAI-87, pp. 68-75,(1987).
- [9] Ait-Kaci, H., Nasr,R.: LOGIN:A Logic Programming Language with Built-in Inheritance: J. Logic Programming, 1,3,pp. 185-215,(1986).
- [10] McSkimin, J.R., Minker, J.: A Predicate Calculus Based Semantic Network for Deductive Searching: Associative Networks, N.V. Findler (eds.), pp. 205-238, Academic Press, (1979).