

記述量圧縮の観点から見た概念体系の構築

松川智義、岸本行生、三池誠司、横田英司、高井貞治、天野真家
(株)日本電子化辞書研究所

概念分類項目間の関係記述データ(項目間記述データ)を、その内容を保存したまま階層化して概念体系を構築する方法を提案する。この方法に基づき概念体系を試作した。その結果、記述量が89%圧縮されること、67%のノードに意味的共通性があることを確認した。これにより、本概念体系が、a. カテゴリ間記述データの検証作業の効率化、b. EDR概念体系における上位概念設定のための予備実験、に用いることができることがわかった。

Construction of Hierarchical Concept Classification Based On Compaction of Concept Description

Tomoyoshi MATSUKAWA, Yukio KISHIMOTO, Seiji HIJKE,
Eiji YOKOTA, Sadaharu TAKAI, Shin-ya AMANO

Japan Electronic Dictionary Research Institute, Ltd.

The method for construction of concept classification by hierarchizing the description data of the relations among concept categories is proposed. The hierarchical concept classification was built by using this method. It was confirmed that the amount of description was compacted at 89.2 %, and that the 67 % nodes in the classification had semantic consistency. With the results, we confirmed that the method can be used for the followings: a. improving efficiency of the verification of the description data of the relations among concept categories, b. pre-experiments for induction of super-concepts in the EDR Concept Classification.

1. はじめに

様々な概念間の関係に関する大規模な知識を信頼性の高い状態で計算機上に実現することは、自然言語処理技術の大きな課題となっている。そのような知識の表現形態としては、概念を階層状に分類整理し、概念間の関係はその階層内に表現するという形態がある。これは、シソーラスと総称される。

田中・仁科は、自然言語処理に必要となるシソーラスとして、上位/下位関係、部分/全体関係を表すシソーラスなどと共に、第3層のシソーラスとして「動詞の世界から名詞の世界を分類したもの」を挙げている^[1]。また、荻野らは、動詞と共起する格要素との関係に着目した動詞の分類を試みた^[2]。このような、他の概念に対する振る舞いに基づく概念の階層的な分類は「概念体系」と呼ばれ、シソーラスの中でも特に重要なものの1つである。

概念体系の構築の際には、次の3点が問題となる。

1. 体系構築のための基礎データの準備、
2. 体系内に記述された知識の効率的な検証、
3. 体系内の節点の導出。

1. としては「実例文の人手による分析データ」が考えられるが、網羅的な概念間の関係を実例文だけから得ることは困難であるので、概念体系構築のための基礎データとしては不十分である。一方、「概念の分類カテゴリ間の関係を内省に基づいて記述したデータ（カテゴリ間記述データ）」の場合、概念間の関係を網羅的に記述することが可能であるので、体系構築のための基礎データとして有効である。ここで、カテゴリ間記述データとは、2つのカテゴリを1つの関係ラベルで結んだもの（カテゴリ間記述タプル）の集合であるとする。

一方、基礎データ内の記述の妥当性を検証する作業においては、個々の記述ごとの検証が不可欠であるので、検証の作業量はデータの記述量に比例して大きくなる。したがって、2. のためには、基礎データの記述内容を保存したまま階層化によって記述量を圧縮し、その上で検証作業を行なうことが有効である。

また、概念体系を節点の階層として構築することの利点の1つとして、節点ごとに概念を分類整理しておくことで、記述の保存を節点単位で行なうことが可能となり、体系の保守管理が容易になるといことが挙げられる。その際に、節点に分類された概念の内容を把握するためには、それらの概念に意味的な共通性があることが望ましい。したがって、3. において導出される節点は、「振る舞い」だけでなく、「意味」的にも共通性がある概念を分類できるものである方がよい。

そこで、本稿では、a. 記述量の圧縮によって記述の効率的な修正を実現し、b. 意味的な把握が可能な概念体系の節点を導出するための、カテゴリ間記述データに基づいた概念体系の構築法を提案する。

以下では、2で、概念体系構築の方法を説明し、3で、その方法を用いてEDRで作成された概念分類項目間の関係記述データ（項目間記述データ）^[5]から概念体系を試作した結果について述べ、4で、作成した概念体系を利用してカテゴリ間記述データを洗練させる方法について議論する。

2. 概念体系構築の方法

本稿で報告する概念体系の構築法は、以下の3つの工程から成る。

- a. カテゴリ間記述のグループ化
- b. グループの階層化
- c. 下位ノードの付加

2. 1~3で、それぞれについて説明する。

2.1 カテゴリ間記述のグループ化

カテゴリ間記述データのような2項関係の実例データは、1/0行列とみなすことができる。このような1/0行列の部分行列で、すべて1で埋まったもの（完全部分行列）やある程度1で埋まったものうち、意味的共通性が認められるものをグループと呼ぶことにする。また、1/0行列内からグループを探し出す操作をグループ化と呼ぶことにする。

```

c#<映像>      <-obj- c#<明度の下降>
c#<映像>      <-obj- c#<明度の上昇>
c#<映像>      <-obj- c#<光の明暗の値>
c#<照明器具> <-obj- c#<明度の下降>
c#<照明器具> <-obj- c#<明度の上昇>
c#<照明器具> <-obj- c#<光の明暗の値>
c#<光>        <-obj- c#<明度の下降>
c#<光>        <-obj- c#<明度の上昇>
c#<光>        <-obj- c#<光の明暗の値>
  
```

(a) カテゴリ間記述データ

G27

```

c#<映像>
c#<照明器具>
c#<光>
  
```

```

c#<明度の下降> 1 1 1
c#<明度の上昇> 1 1 1
c#<光の明暗の値> 1 1 1
  
```

(b) カテゴリ間記述のグループ化

G27<-obj-G27

```

列側
G27 : (c#<映像> c#<照明器具> c#<光>)
行側
G27 : (c#<明度の下降> c#<明度の上昇> c#<光の明暗の値>)
(ただし、「A:(c#B ...)」は、「A<-instance_of-c#B」の略記)
  
```

(c) 所属関係による項目間記述の階層化

図1. カテゴリ間記述のグループ化

カテゴリ間記述データをグループ化した後、グループの列側、行側にそれぞれ記号を与えれば、グ

ループに属するカテゴリ間記述をそれらの記号間の記述(グループ間記述)として抽象化して表現することができる。これは、グループに属するカテゴリの上位概念を設定し、上位概念と各カテゴリの間に所属関係(instance_of で表現)をつけて階層化することに対応している(図1)。

このとき、グループの行側/列側に属するカテゴリの個数をそれぞれk, lとすると、

$$k l - (k + l + 1)$$

だけ、記述量が減少する。カテゴリ間記述ダブルがk lだけ減るかわりに、グループの所属関係記述ダブルが(k+1)、グループ間記述ダブルが1、それぞれ増えるからである。

そこで、記述量圧縮のために、カテゴリ間記述データをグループ化する。グループ化は、a. グループの自動抽出、b. グループの修正、という手順で行なう。以下で、a. b. それぞれについて説明する。

2. 1. 1 グループの自動抽出

単語共起データからボトムアップ的に単語のグループ化を行なう方法として、筆者の一人はDM分解算法を利用した方法を提案した^[7]。この方法は、一般の2項関係データのグループ化に応用することができる。そこで、この方法を用いてカテゴリ間記述データのグループ化を行なう。

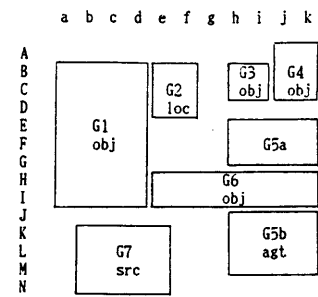
ただし、カテゴリ間記述についてのグループは、完全部分行列であるとする。これは、文献[7]で単語共起についてのグループを、ある程度1で埋まった部分行列としていたことと異なる。単語共起データの場合、可能な単語共起を網羅したデータを得ることは不可能である。これに対してカテゴリ間記述データの場合は、網羅的にカテゴリ間記述を行なうことが前提となっている。そこで、カテゴリ間記述データに対しては、完全部分行列を探すことでグループ化を行ない、元のデータの持つ情報を保存することにした。なお、完全部分行列を求めるためには、文献[7]のアルゴリズムにおいて、パラメータ「1-0反転率q」を0と設定すればよい。

また、文献[7]では、完全部分行列(意味素性と呼んでいた)を大きさ(グループに属する要素の個数)の大きいものから探していた。単語共起データを表現したモデル上での動作を保証するためであった。しかし、カテゴリ間記述データの場合はこの方針にはこだわらず、アルゴリズムの実行の途中で大きい値(例えば6)よりも大きな完全部分行列が見つかることに出力することにする。

2項関係データからグループを自動抽出した結果の例を図2(a)に示す。これから、図2(b)(c)の抽象化した記述が得られる。

2. 1. 2 グループの修正

グループの修正は、a. グループにカテゴリを付加したり、b. グループからカテゴリを削除した



(G5a, G5b を含ませて、G5 であることに注意)

(a) グループ化の結果

列側	行側
G1 : (a b c d)	G1 : (B C D E F G H I)
G2 : (e f)	G2 : (B C D)
G3 : (h i)	G3 : (B C)
G4 : (j k)	G4 : (A B C)
G5 : (h i j k)	G5 : (E F G J K L M)
G6 : (e f g h i j k)	G6 : (H I)
G7 : (b c d e)	G7 : (K L M N)

(b) 所属関係の記述

G1 <-obj- G1
G2 <-loc- G2
G3 <-obj- G3
G4 <-obj- G4
G5 <-agt- G5
G6 <-obj- G6
G7 <-src- G7

(c) グループ間の関係記述

図2. グループ化に基づく抽象化した記述

りすることを繰り返すことで行なう。

カテゴリの付加

グループに含まれていないカテゴリで、意味的な共通性と振る舞いの共通性から考えて、グループに含まれていてもよいと思われるものがあれば、そのカテゴリをグループに付加する。

カテゴリの削除

グループの中に、他のカテゴリと比較して意味的な共通性の認められないカテゴリが含まれている場合、あるいは、カテゴリ間記述データの不備から振る舞いの共通しないカテゴリが含まれている場合は、そのカテゴリをグループから削除する。

2. 2 グループの階層化

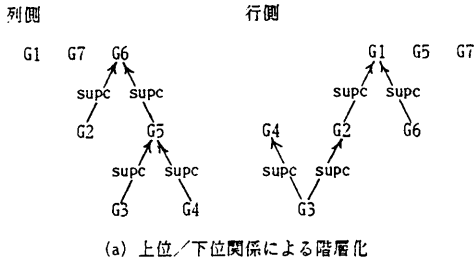
図2(b)にも示したように、カテゴリ間記述をグループ間記述に置き換える場合は、カテゴリのグループへの所属関係を別に記述する。この所属関係の記述量を、グループ間の包含関係(上位/下位関係)の記述を導入することによって減少させることができる。具体的には、以下のようにする。

1. 列側、行側それぞれについて、グループ間の包含関係を調べる。
2. グループを節点とし、グループ間の包含の直前/直後関係を枝とする半順序グラフを作成する。

その枝ごとに上位/下位関係 (supc で表現) をグループ間の関係として記述する。以降、グラフの節点としてのグループをノードと呼ぶ。

- 各ノードについて、ノードに属するカテゴリで下位のノードにも重複して属しているものは、その所属関係の記述を削除する。

各グループに属するカテゴリは、ノードの上位/下位関係を下位 (子孫) に辿っていくことによって得ることができる (図3)。



列側	行側
G1 : (a b c d)	G1 : (X Y Z E F G X X)
G2 : (e f)	G2 : (X Y D)
G3 : (h i)	G3 : (B C)
G4 : (j k)	G4 : (A Y Z)
G5 : (X Y Z X)	G5 : (E F G J K L M)
G6 : (Y Z s X Y Z)	G6 : (H I)
G7 : (b c d e)	G7 : (K L M N)

(b) 所属関係の記述の削除

図3. グループの階層化

上位/下位関係にある2つのノードの集合差の大きさを m とすると、3.において、上位/下位関係の記述1つごとに、カテゴリのノードへの所属関係の記述が m 個減る。グループ化の方法から、 $m \geq 1$ が言えるから、グループの階層化によって、全体の記述量は減少こそすれ増加することはない。

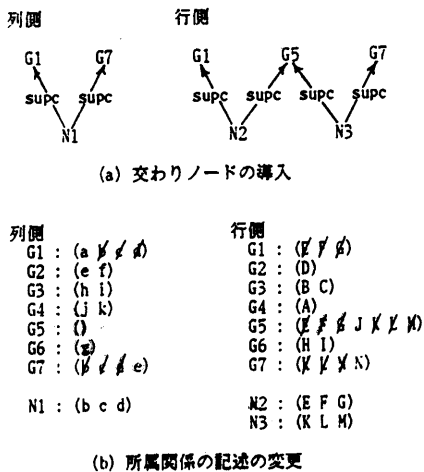


図4. 下位ノードの付加

2.3 下位ノードの付加

交わりのある2つのノードの下位に新たにノードを導入することによって、さらに所属関係の記述量を減少させることができる (図4)。具体的には、以下のようにする。

- 交わりが2以上ある2つのノードを見つけて、それらの下位に新たなノード (交わりノード) を導入する (supc で表現)。
 - 交わりに属するカテゴリを交わりノードに所属させる (instance_of で表現) と同時に、上位の2つのノードにおけるそれらの所属関係の記述を削除する。
 - 交わりノード間の包含関係を調べて、上位/下位関係で表される部分を修正する。
- 1.で、交わりを2以上としたのは、2つの上位ノードと交わりノードとの間の上位/下位関係のための記述量が2だけ増加するのを打ち消すためである。交わりに属するカテゴリの個数を n とすると、交わりノードを導入するごとに、 $(n-2)$ だけ所属関係の記述量が減少する。

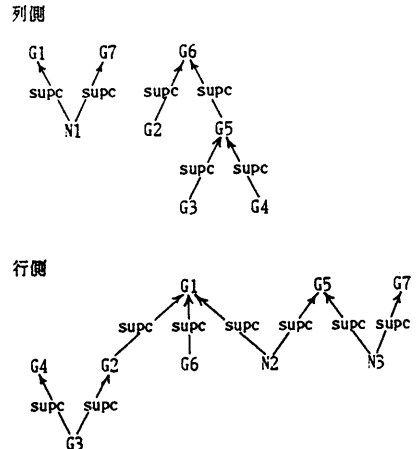


図5. 概念体系

この作業では、1. における交わりのある2つのノードの選び方に任意性がある。したがって、交わりが2以上あるノードの対をすべて計算して求めた後、その中から交わりノードの選択を手で行なうという方法も考えられる。

また、交わりノードをこの作業の対象に随時加えていくことにより、さらに記述量を減少させていくことも考えられる。しかし、明らかに計算量の爆発を引き起こし現実的ではないので、これは行なわない。

以上の工程により概念体系を得る。図5に、図2(a)のデータから得られる概念体系を示す。なお本体系の作成は、2項関係抽象化支援システム^[6]を用いて、計算機支援の下に行なうことが可能である。

3. 概念体系の試作結果

EDRでは、概念の分類カテゴリである概念分類項目(項目)を、モノ概念について約200個、コト概念について約250個、それぞれ設定し、それらに対して約30万の概念を分類する作業を行った^[4]。一方、可能な項目間の関係を人手で記述して、網羅

的な項目間記述データを作成した^[5]。

今回、この項目間記述データを対象に概念体系の試作を行なった。その結果、39,256タプルの項目間記述データのうち、36,562タプルがグループ化され、2,694タプルがグループ化されずに残った。グループの個数(ノード間記述の記述数)は75、交わりノードの個数は48(列側37、行側11)、上位/下位関係の記述数は178(列側127、行側51)、所属関係の記述数は1,279(列側978、行側301)であった。したがって、記述量の圧縮率は、

$$(39,256 - (75 + 178 + 1,279 + 2,694)) / 39,256 = 89.2\%$$

になる。

試作した概念体系を図6に示す。列側、行側それぞれの上位/下位関係、及び、ノード間の関係記述をネットワーク状に表現し、項目のノードへの所属関係の概略をノードの横に付記した。また、一部のノードには、ノードの子孫を総称する名称を()で囲んで付記した。

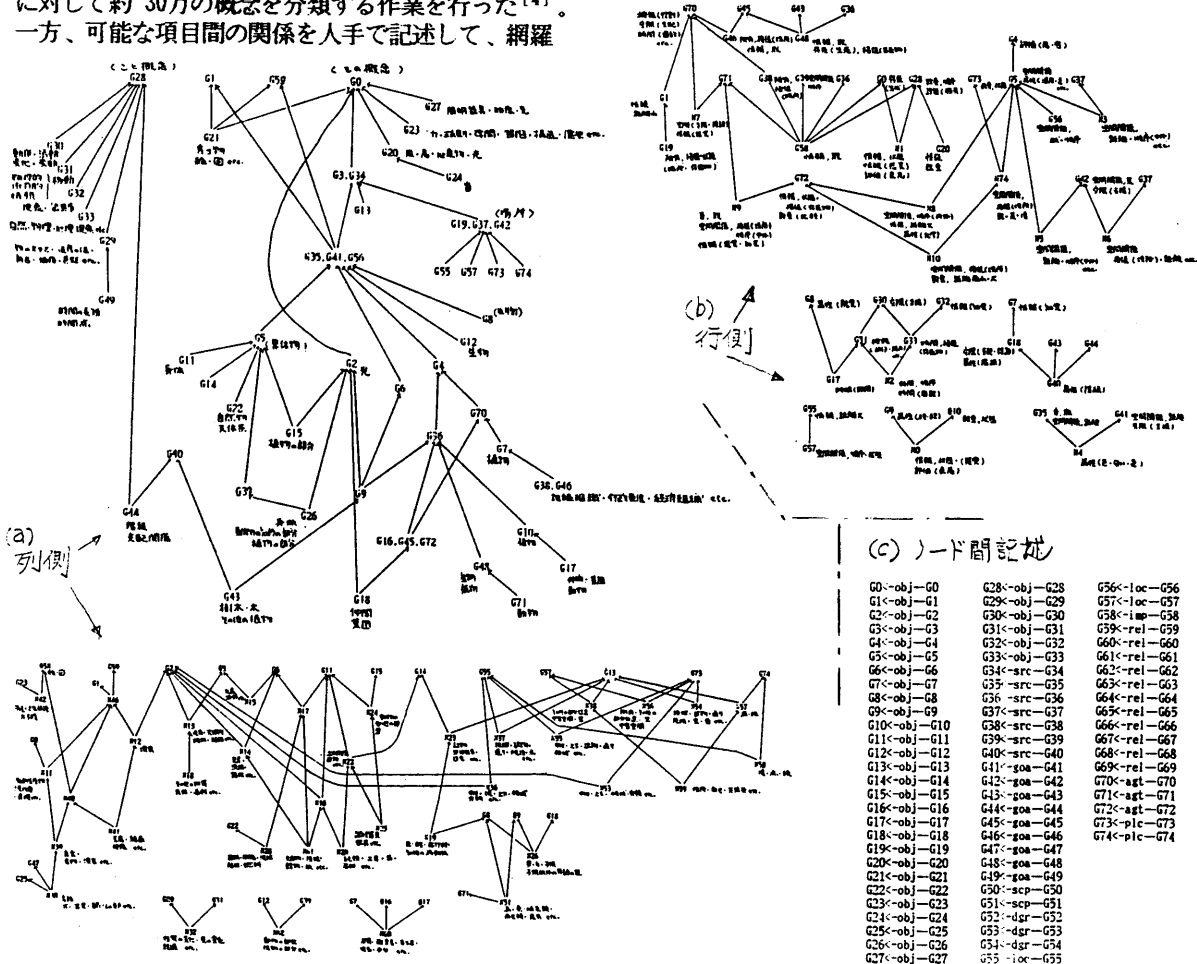


図6. 概念体系の試作結果

本概念体系におけるノードは、項目間記述データにおいて振る舞いの共通する項目がまとめられたものである。しかし、1でも述べたように、体系の保守管理及び項目間記述データ検証の観点からは、ノードに属する項目の間に、意味的な共通性が存在することが望ましい。そこで、試作した概念体系内の各ノードに属する項目に意味的共通性があるかどうかを調査した。その際、意味的共通性は、例えば、

N28: (c#<固体>、c#<樹脂>、c#<液体>、
c#<粘体>、c#<燃料>)

に対する「物質」や、

N2 : (c#<時間的順序が変化する>、
c#<回数>、c#<回数が増える>、
c#<ある時間にコトが満ちる>)

に対する「時間」のように、項目を総称する簡潔な表現があるかどうかで判断した。結果を表1に示す。

表1. 意味的共通性によるノードの分類

	列側ノード	行側ノード
意味的共通性あり	57	39
意味的共通性なし	13	34
所属項目数 1	10	8
所属項目数 0	32	5

所属する項目数が0及び1のノードを計算から除いたとしても、半数以上(67%)のノードに意味的共通性が認められたことがわかる。特に、列側のノードは約8割のノードに意味的共通性が認められた。したがって、本体系を1で述べた目的に用いることが可能であると言える。

4. 本概念体系の利用法

1でも述べたように、本概念体系構築の主な目的は、a. カテゴリ間記述データの効率的な検証、b. 概念体系内の節点導出のための予備実験、である。以下では、特にa. に関して前節の、項目間記述データを基にした概念体系の試作結果(試作体系)を例にして説明する。

本概念体系を利用することにより、項目間記述データの効率的な検証が可能になる主な原因は以下の5点である。

- ノード単位で検証できること。
- 適切なレベルにおける検証が可能なこと。
- 体系における項目の相対的な配置に基づく検証(大域的な検証)が可能なこと。
- 階層外記述に対する検証が行なえること。
- 体系構築作業を分業化できること。

以下に、それぞれについて詳しく述べる。

ノード単位の検証

本体系においては、項目間の関係はすべてノード間記述として記述されるので、項目間記述データの検証もノードごとに行なうことが可能である。ノード間記述の表現している項目間記述の集合は、ノードの子孫を辿ることにより図7のように得られる。すなわち、図7のようにノード間記述の行側・列側のノードをそれぞれ子孫全体について展開して得られるすべての項目の組み合わせについて、元のノード間記述と同じ関係が成立する。ノード単位で項目間記述を検証する場合、このように子孫の項目を展開して一覧した上で、検証を行なう。

G5<obj-G5

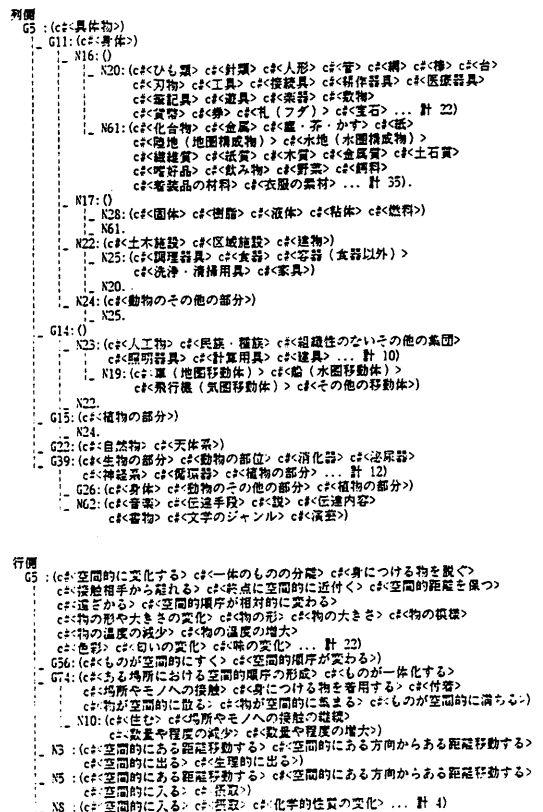


図7 G5の子孫

3で示した通り、本試作体系の半数以上のノードには意味的共通性がある。また、意味的共通性がないノードに関しても、ノード内をいくつか分割すれば、そのそれぞれについて意味的共通性が認められる場合が多い。したがって、各ノードの子孫全体を1つのものとして把握することが比較的容易である。

例えば、図7. において、列側は、G5 以下は「具体物」、N20 以下は「形状、工具・器具、券など」、N61 以下は「材質、陸・水地、飲食物、材料など」、N23 以下は「集団、器具、移動体」、G39

以下は「生物の部分、音楽・書物・演芸など」と把握できる。一方、行側は、G5 以下は「空間関係」、G74 以下は「順序形成、接触、集散、満ちること、など」、N3 以下は「移動、出ること、など」と把握できる。(ここでは便宜的に把握の仕方をして「」で囲って表現した。しかし、実際にノードの子孫を1つのものとして把握する際には、ノードの子孫に属する項目に関する具体的な印象があれば十分であり、特に言葉にする必要はない。)

ノードをそれぞれ1つのものとして把握してノード間記述タプルを検証する。その場合、項目間記述タプルをランダムに検証する場合と比べて、記述量の圧縮率だけ作業の効率化が望める。ノード間記述タプルごとに、その子孫の項目間記述タプルの正否について一括して検証できるからである。例えば、図7のノード間記述タプル

G5 <-obj- G5

の検証を、列側、行側のノードG5をそれぞれ、「具体物」「空間関係」として把握し、

「具体物」<-obj-「空間関係」

の妥当性を検証することで、このノード間記述の子孫の3,990個の項目間記述タプルを一括して検証できる。

適切なレベルにおける検証

前節のようなノード単位の検証だけでは、ノード間記述の十分な検証が行えない場合がある。ノードの把握の仕方が粗すぎると感じられる場合がある。その場合、子孫ノードのレベルでさらに細かい検証を加えることが必要になる。このときも、子孫ノードをそれぞれ1つのものとして把握することができれば、やはり作業の効率化が望める。例えば、先のノード間記述

G5 <-obj- G5

の妥当性を検証する場合、

G11 G5(ノードのみ)

G14 <-obj- G56, 74

G15, 39 N3, 5, 8

のように、子ノード間の記述に展開してから、各子ノードを1つのものとして把握した上で、それぞれの記述を一括して検証する。

また、子ノード間記述を検証する代わりに、ノードの子孫の項目を流し読みして、ノードに対して抱いていた印象と相違がないことを確認するという方法も考えられる。このとき、印象と相違する項目に対しては、さらに詳しい検証を加える。

以上のように、項目間記述の検証をレベルに合わせて段階的に行なえるので、ランダムに検証する場合に比べて、a. 無駄な作業を自然に避けられる、と同時に、b. 作業に連続性が生じるので作業者の負担感も軽減される。

大域的な検証

項目間記述タプルの作成作業は、項目の組み合わせごとに関係の可能性を内省することによって行なわれる。その際、他の項目間記述タプルとの整合性は基本的には考慮されない。つまり、項目間記述タプルは、そのタプルに対する局所的な内省のみによって作成される。したがって、概念体系内における各項目の相対的な配置に基づいて、大域的に項目間記述を検証することが必要となる。

反対概念を表すコト項目などは、あらかじめ同じ振る舞いをすると予想される。例えば、c#<認知対象との認知的距離減少>、c#<認知対象との認知的距離増大>、という2つの項目は、それぞれ反対概念として設定されたものである。振る舞いは同じになると考えられる。ところが、図6の概念体系においては、前者は、G1、G20、G23、後者は、G55、N8、とまったく異なるノードに所属している。この差異が妥当であるかどうかは、各項目に対する記述を調べることにより判断しなければならない。

```

= c#<認知対象との認知的距離減少>の記述先 **
(c#<認知対象との認知的距離増大>との差)
- G1の根の記述先 (G4, G5, G72 以下は共通)
G1<-obj- G1
G1:()
  G21:(c#<書き物> c#<札> c#<券> c#<書物> c#<絵> ... 計 8)
  G35: G41, G56:()
  G6:()
    N14:(c#<元菜> c#<生理的移動物質> c#<液体> c#<固体>)
    N6:(c#<化合物> c#<繊維質> c#<嗜好品> c#<飼料> c#<糖> ... 計 35)
  N17:()
  G8:()
  G12:(c#<生物> c#<生物の部分> c#<消化器> c#<臓物の部分> ... 計 13)
G70<-ast- G70
G70:()
  G28, G46:(c#<行政単位> c#<公共性の高いサービス機関> ... 計 9)
  N60:(c#<職業名> c#<姓名> c#<その他の種類からみた人間> ... 計 19)
- G20の根の記述先
G20<-obj- G20
G20:(c#<その他> c#<その他のもの概念> c#<風> c#<心象物> c#<光>)
  G24:(c#<音> c#<人の音声> c#<人の音声以外の音>)
G28<-obj- G28
G28:(c#<性状> 性状) c#<人の社会的属性値> c#<性格> c#<状態の形質> ... 計 51)
  G29:(c#<物の大きさ> c#<具体物に対する感質的属性値> ... 計 35)
  N2:(c#<時間的産物> c#<物質的産物>)
  N3:(c#<容性物質の質化> c#<性質の変化> c#<量の質化> ... 計 66)
  G30:(c#<動作> 活動) c#<感情活動> c#<変化> c#<活動> ... 計 66)
  G31:(c#<移動> c#<所有移動物> c#<情報移動物> ... 計 20)
  G32:(c#<実体現象> c#<物理現象> c#<現象現象> c#<現象現象> ... 計 31)
  G33:(c#<現象現象> c#<心理現象> c#<社会現象> ... 計 15)
  G44:(c#<情報> c#<支配関係> c#<認知していない何か>)
- G23の根の記述先
G23<-obj- G23
G23:(c#<力> c#<法則> c#<空間> c#<習俗> c#<構造> c#<産業> ... 計 28)
  N2:(c#<構造> c#<法則現象> c#<方法手段>)
= c#<認知対象との認知的距離増大>の記述先 **
(c#<認知対象との認知的距離減少>との差)
- G55の根の記述先
G55<-loc- G55
G55:()
  N36:(c#<中心> c#<境> c#<下> c#<側面> c#<穴道> ... 計 11)
  N37:(c#<行政区域> c#<任意地> c#<道> c#<宛先> c#<道> ... 計 23)
  N35:(c#<中心> c#<行政区域> c#<道> c#<領域> c#<区画> ... 計 31)
  N38:(c#<境> c#<点> c#<線>)
- N8の根の記述先 (G4, G5, G72 以下は共通)

```

図8. 反対概念を表す2つの項目に対する記述

各項目に対する記述は、項目の属するノードの祖先を辿っていくことによって得られる。例えば図8では、G1、G20、G23、及び、G55、N8、の祖先ノードの記述を展開することにより、上記の2つの項目に対するすべての記述を得ている。これを基にノード単位の記述の検証を行ない、両項目の振る舞いの

