

国語辞書の意味分類を利用した概念ベースにおける多義概念の分割

山西 公一郎 小島 一秀 渡部 広一 河岡 司

同志社大学大学院 工学研究科 知識工学専攻
〒610-0394 京都府京田辺市多々羅都谷 1-3

本稿では常識判断メカニズムにおいて中核をなす「概念ベース」の多義性の解消について提案する。常識判断メカニズムは“人間らしい常識的な判断や推測”をする知能ロボットへの応用を目的としている。概念ベースとは概念をその他の概念（属性と呼ぶ）集合で表した知識ベースであり、電子化辞書等から機械的に構築される。このため雑音とよばれる不適切な属性も多く含まれる。また多義性を考慮して構築されていないための問題もある。本研究の目的は国語辞書の意味分類をもとに概念ベース内の多義概念の属性を意味によって分割し、新しい多義性のない概念ベースに改良することである。提案方式によって作成した多義性のない概念ベースはサンプルによる目視判断において改善されていることを確認できた。また関連の近さを示す評価尺度用テストデータを用いた実験においても改善されていることを示した。

Division of the concept which has ambiguity with the meaning classification of language dictionary

Koichiro Yamanishi, Kazuhide Kojima, Hirokazu Watabe and Tsukasa Kawaoka

Graduate School of Engineering, Doshisha University
Kyotanabe, Kyoto 610-0394

Human can judge intelligently with imperfect information than computer. A concept-base is one of the main elements to realize the intelligent judgement by computer. This paper shows that the concept-base becomes closer to human judgement by division of the polysemous concept. The concept-base is a knowledge base in which each concept consists of a set of concepts (attributes). Concretely with the meaning classification of language dictionary, each attribute is divided to proper meaning category. It is shown that the concept-base after the division is more effective for the common sense judgement system by the experimental result using degree of association.

1. はじめに

情報処理技術の発展は目覚ましいが人間らしい知的な処理の実現にはまだまだ多くの問題が残されている。そこで、人間的な常識判断をコンピュータが行えるようにするというのが本研究の目的である。常識判断には、物の大小、長さ、広さ、重さ、場所、時間、速さとい

った量的な判断、また赤い、熱いといった感覚判断、そして、うれしい、悲しいといった感情判断などがある。これら常識判断を実現するのが常識判断メカニズムである。コンピュータが人間のように柔軟に判断するには人間の持つような語に関するさまざま知識、語の“概念”を与える必要がある。しかし、実世界の概念の

数は膨大であり、人間が手作業で概念ベースを構築することは大変な作業である。このため、概念ベースを機械的に構築する必要がある。

概念ベースにおいて、概念はそれを説明する属性と呼ばれる概念とその重みで定義される。このように構築された概念ベースにおいては、多義に関する考慮はされていない。しかしながら人間が扱う言葉の概念はさまざまな観点、ニュアンスを含んでいる。コンピュータにこの多義の概念を理解させるのは不可能であるため、概念を意味で分割して概念の持つ多義性を消す必要がある。次の二つの文章を考える。

文1：星がまたたいている。

文2：星をつかまえる。

この2文の“星”は意味が違う。文1の“星”は天体の意味であるが、文2の“星”は犯人という意味である。多義性による分割をしなければ、知的ロボットに文2の文章を伝えたならば、空の星をつかまえようとする??といった行動を起こしてしまう。本稿では常識判断メカニズムにおいて上のような入力があった際に概念の違いを区別することができるような概念ベースを作るための多義概念の分離法を提案する。

2. 概念ベース

2.1 概念ベースと常識判断メカニズム

概念ベースは常識判断メカニズムで用いられるが、ここでは常識判断メカニズムにおける概念ベースの位置付けについて述べる。

常識判断メカニズムにおいて概念ベースは中核をなすものである。図1に本研究の目的である常識判断メカニズムの全体図を示す。図1を左のユーザ側から見ていくと、まずユーザと直接対話を行う会話メカニズムがある。会話メカニズムは背後にある様々な判断メカニズムを用いて人間との会話を行う。判断メカニズムはそれぞれ独自の処理方式や知識ベースを持っていてそれぞれの判断を行う。判断メカニズムは入力が自分の知識ベースにないときなどは、その背後にある連想メカニズムから概念に関する知識を得て処理を行う。連想メカニズムはその背後の概念を定義する概念ベースを用いて連想処理を行う。

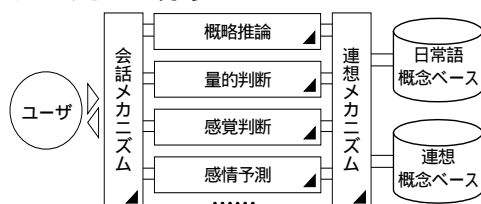


図1 常識判断メカニズムの全体像

2.2 概念ベースの多義性

概念ベースは複数の国語辞書から自動構築された知識ベースである^[1]。概念はその概念を説明する複数の概念(属性と呼ぶ)で定義している。そして属性には出現頻度による重みがついている。概念数は約3万、属性数は約150万である。

概念ベースは辞書の意味分類を無視して作成されているため、見出し語概念が複数の意味を持っていたとしても属性側ではその違いを表すことができていない。すなわち、概念ベースは多義概念に対する属性の区別が全く考慮されていない。多義概念とは複数の意味を持つ概念のことである。複数の意味とは、比喩的に言う意味も含んでいる。概念ベースの雑音には完全な雑音とこの多義概念のための雑音がある。

多義概念のための雑音とは、意味が複数あるために一方の意味からは雑音と判断される属性ともう一方の意味からは適切な属性となるような属性のことである。多義概念の雑音に関しても概念を意味によって分割して属性を分けることで、より正確な意味判断が可能となる。

多義概念の例を下に示す。例は星の属性であるが、重みが小さいほど雑音である可能性は確かに高くなっている。最も低い属性“歌”は星の属性としては不適切である雑音と考えられる。しかし、重みが高い“鋌”も明らかに雑音でまた、“犯人”は星の“天体”という意味からすると雑音である。しかし、“犯人”も“星”が持つ“容疑者”の意味からすると適切な属性であり、逆にそのとき“天体”が雑音となる。“天体”、“犯人”、どちらも星の属性としては非常に適切な属性ではあるが、区別せずに格納することで問題が生じてしまう。これが多義性雑音である。この問題は意味により属性を分離することで解決する。

星 = {(星, 122), (天体, 78), (恒星, 59), (惑星, 53), (点, 51), (鋌, 49), (犯人, 48), (勝ち負け, 44), (容疑者, 43), ..., (歌, 14)}

概念ベースの概念数は33,699であるが、そのうち、3.1節で述べる概念構造情報によって多義概念であるとされる概念は全部で10,603個あった。意味数と概念数の関係を図2に示す。

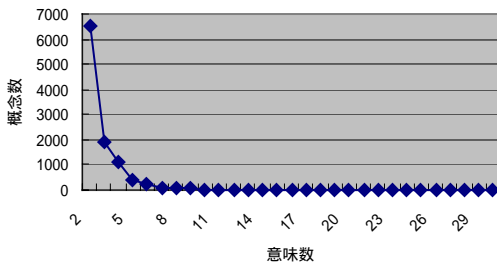


図2 意味数による概念数

2.3 概念間の関連度計算

本研究では概念分割や概念ベースの評価に関連度を用いている。関連度とは、概念間の関連の強さを定量化した値であり、概念の属性(1次属性)と属性の属性(2次属性)の一致度により求める^{[2][3]}。以下では、一致度に基づく関連度の計算方法について述べる。

一致度は概念の1次属性がどの程度一致しているかを示す0から1の値で、以下のように計算する。一致度を求める2概念をA, Bとする。概念A, Bは下式の通りである。

$$A = \{(a_1, u_1), (a_2, u_2), \dots, (a_L, u_L)\}$$

$$B = \{(b_1, v_1), (b_2, v_2), \dots, (b_M, v_M)\}$$

ただし, L, MはA, Bの属性数である。このときの概念A, Bの一致度 match(A,B)は次のようになる。

$$\text{match}(A,B) = \frac{1}{2} \left(\frac{u_M}{U} + \frac{v_M}{V} \right)$$

ただし, u_M は $a_i=b_j$ が存在する u_i の合計, v_M も同様で $a_i=b_j$ が存在する v_j の合計である。U, Vは次のようになる。

$$U = \sum_{i=1}^L u_i$$

$$V = \sum_{i=1}^M v_i$$

関連度は概念間の関連の深さを示す0から1の値で、以下のように計算する。関連度を求める2概念をA, Bとする。属性数の多い方をAとする。したがって, $L > M$ である。

まず, A, Bの属性を一対一で対応付ける。このとき, 対応付けられた属性の一致度の合計が最大になるようにする。これは組み合わせ最適化問題とつながるが実際には計算量の関係

から、属性の全ての組の中で最も一致度の大きな組から対応を決めている。Bはそのままにし、Aの属性を並べ替えてA'を作る。A'の第i属性はBの第i属性と対応する。Bの属性に対応しなかったAの属性は無視する。したがって、A'の属性数はMとなる。

$$A' = \{(a'_1, u'_1), (a'_2, u'_2), \dots, (a'_M, u'_M)\}$$

関連度 Rel(A,B)は次のようになる。

$$\text{Rel}(A,B) = \frac{1}{2} \left(\frac{u'_m}{U} + \frac{v_m}{V} \right)$$

ただし, u'_m, v_m は次のようになっている。

$$u'_m = \sum_{i=1}^M \text{match}(a'_i, b_i) u'_i$$

$$v_m = \sum_{i=1}^M \text{match}(a'_i, b_i) v_i$$

3. 概念分割

概念分割とは現在の多義性のある概念ベースから多義性を解消し、1概念1意味の概念ベースを作ることである。多義の概念ベースの属性を意味により分割し(属性の意味分割)、そして分割された属性それぞれに概念を付与しなおす。具体的には概念の表記は同じであるが異なる概念の識別詞を付与する。(属性の意味決定)。Ex.星(1001):天体。星(1002):犯人。

3.1 概念構造情報

今回の概念ベースの分割には概念構造情報を用いた。これは概念ベースと同様電子化辞書から機械構築したものであるが概念ベースとは異なり、意味毎の詳細な分離情報を保存している。^[4]

概念ベースは基本的に辞書の一つの見出し語が一概念となっているのに対して、概念構造情報の一つは構築に用いた辞書の一つの分類に対応する。今後概念ベースの概念と区別するため、概念構造情報の概念を見出し語と呼ぶ。また概念ベースの属性と区別するため、概念構造情報では関連語と呼ぶ。国語辞書では見出し語の説明にある記号や表記特徴を利用して、見出し語との関係が分かる。概念ベースとの大きな違いは見出し語とそれを説明する

関連語とのあいだに明確な論理関係が記載されていることである。論理関係の種類は、同義、類義、上位、対義、尊敬、丁寧、英字である。また国語辞書は見出し語を説明する欄では意味の違いによって番号が付けられている。それを利用して概念構造情報には自動構築する際に各関連語に意味番号をつけている。つまり、見出し語が多義語の場合、関連語は見出し語のどの意味の関連語であるかが分かるのである。見出し語に対して今回は概念構造情報の中でも特にこの意味番号を利用することにした。

見出し語数は約 16 万、関連語数は約 100 万、1 見出し語あたりの平均関連語数は 6、一つの見出し語が持つ関連語数は 0 から 97 となっている。概念構造情報の例を表 1 に示す。

表 1 概念構造情報の例 “星”(見出し語)

意味番号	関連語(関連語, 関係型)
1	(星雲, 類義)(天体, 上位)・・・
2	(階級, 不明)(記憶, 不明)・・・
3	(点, 上位)(斑点, 上位)・・・
4	(眼球, 不明)(白い, 不明)・・・
5	(成績, 上位)(白星, 上位)・・・
6	(目当て, 同義)(目ぼし, 同義)・・・
7	(犯人, 同義)・・・
8	(運勢, 同義)(生まれる, 不明)・・・
9	(移る, 不明)(形, 不明)・・・
10	(スター, 同義)(花形, 同義)・・・

3.2 概念ベースと概念構造情報の対応

分割処理では初めに概念ベース(概念)と概念構造情報(見出し語)との対応を取る。

見出し語と概念の対応だが、見出し語の対応とは、概念ベースの概念名“星”に対応するのは概念構造情報の“星”である。概念数は約 3 万概念である。一方概念構造情報の見出し語は約 9 万語である。対応を次のように決定する。概念ベースの概念や、概念構造情報の見出し語はそれぞれ表記を複数持つ。例えば、概念“星”の場合は(ほし)と(星)である。また概念“木”の表記は(木)と(き)である。対応を取る際、表記の(き)からは概念“木”以外に“気”も考えられる。そこで対応の取り方のアルゴリズムを以下のようにする。

概念ベースの概念 A の概念名の表記を a_1, a_2, \dots とし、概念構造情報の見出し語 S_1 の表記を $s_{11}, s_{12}, \dots, s_{21}, s_{22}, \dots$ とする。 a_1 と一致する表記が s_{11}, s_{41} とすると候補見出し語は S_1, S_4 となる。また a_2 のほうも同様にして候補見出し語が s_{23}, s_{42} だったとする。このとき概念 A に対応する候補となる

概念構造情報の見出し語は S_1, S_2, S_4 であるが表記の候補となった見出し語で S_4 は二度候補にあがっている。従って概念 A の対応する概念構造情報の見出し語は S_4 というようになる。図 3 に概念対応の例を示した。以上のようにして、概念に対応する見出し語を決める。対応の結果を表 2 に示す。

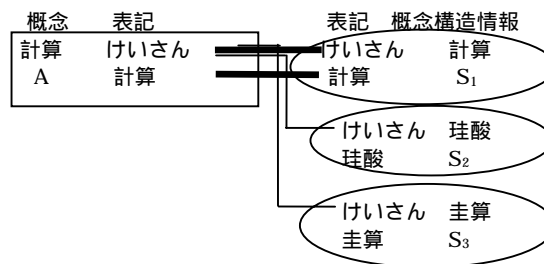


図 3 概念ベースと概念構造情報の対応の取り方

表 2 概念ベースと概念構造情報の対応の結果

対応の取れた概念数	31,244
対応の取れなかった概念数	2,455
一意に決まらなかった概念数	1,415
(計)全概念数	33,699

表 2 の対応が取れなかった概念とは対応する概念構造情報(見出し語)がない場合である。この概念に関しては機械的な分割は不可能となる。また、一意に決まらなかった概念は、一つの概念に対して複数の概念構造情報の見出し語が存在する場合である。例えば概念“円”である。概念“円”の表記に(円, えん, まどか)とあり、概念構造情報の見出し語(円, えん)と(円, まどか)二つに対応が取れてしまう。これは概念ベース作成時に簡単な統合が行われているためである。この場合、概念構造情報側で(円, えん)の意味数が 3 個、(円, まどか)の意味数が 2 あるのである。対応する概念構造情報の見出し語を複数にする。その場合、意味数は $3+2=5$ となる。

3.3 1 次属性の意味分割方法

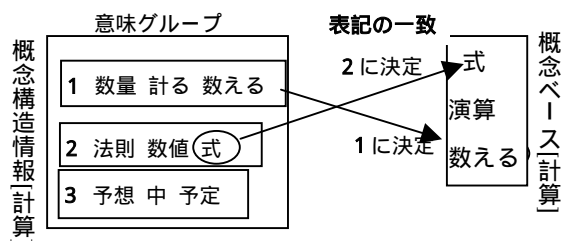
概念ベースの概念と概念構造情報の見出し語の対応が取れたら、概念 A の属性をそれに対応する概念構造情報の見出し語 S をもとに分割する。分割では表記の比較、関連度計算という順に処理を行う。

1) 概念 A “計算” の属性と見出し語 S の関連語を比較し、関連語と同じ表記がある属性

に関してはその関連語の意味番号をふる(図4).

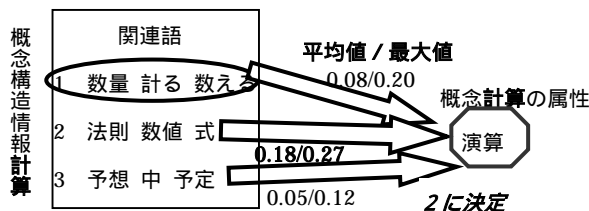
- 2) 1) で意味番号が決まらなかった属性と全関連語と関連度を計算し, 各意味番号の関連度^[3]を計算し一番適切な意味番号に決める(図5).
 - 2) の一番適切な意味番号は, 各属性と意味番号ごとに値を計算し, 最大となった意味番号が一番適切な意味番号とする. 意味ごとの値は最大値と平均値の二方法をとった.
- 例として概念“計算”を分割する.

計算={ (式,37), (演算,36), (数える,11), (演算,11), ..., (原価,6) }



属性“式”と“数える”は構造情報の関連語と表記が一致するので意味番号がこの時点で決まる

図4 ステップ1 概念計算の表記での分割



“演算”の意味番号を決める
 平均値分割では2番目に決定
 最大値分割でも2番目に決定

図5 ステップ2 概念“計算”の関連度での分割

1次属性分割結果は以下ようになった.

表3 概念計算の分割結果

意味番号	意味	属性
1	数量を計る	算盤 算定 算珠 算筆 算複 利考慮 算用 計算 逆算
2	数式	演算 式 計算 検算 精算 運算 算木 通計 数理 合算
3	予測	結果 予測 胸算用 算入 誤算 出す 暗算 見積もる

3.4 属性の意味決定

次に各1次属性が2次属性(多義性を持つ)の何番目の意味に対応するかを決める(図6). 2次属性とは概念の1次属性を概念と見た属性のことで, 基準の概念から見ると2次属性ということになる. 星の1次属性“太陽”の意味グループが決定されるので2次属性についても“太陽”の意味グループに対応する“夕日”の意味グループを抽出する必要がある. これによって分割された概念の一義性が実現する. この際, 1次属性が多義概念でない場合は処理の必要がない.

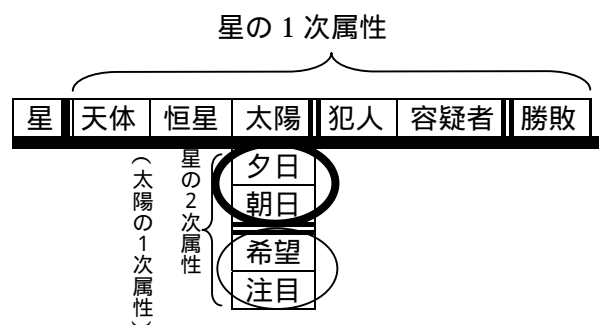


図6 1次属性の意味決定

属性の意味決定もまず表記の一致で意味番号を決め, 決まらなかった属性に対して, 関連度分割での決定を行う.

表記での決定方法は, 意味を決める1次属性 a_i の属する意味番号を X とすると a_i に属する他の1次属性を対象にする. そして a_i の属性つまり2次属性すべてと X に属する1次属性とを比較して表記が一致したものととする. そのときその一致した2次属性が属する意味番号が a_i の持つ属性の意味とする.

概念“計算”を例に説明する.

- 表記での決定
 計算の意味1の1次属性は以下となる.
 計算1 = { 算盤, 算定, 算 }

算盤は2つの意味を持っている.
 算盤1 = { 算, 盤面, 位置 }
 算盤2 = { 損得, 発明, 利益 }

算定の意味は一つだけである.
 算定 = { 数える, 決まる, 計画 }

算は次の6つの意味を持つ.

- 算 1 = { 木片, 用具, 四則 }
- 算 2 = { 歳, 式, 演算 }
- 算 3 = { 算木 }
- 算 4 = { 数える, 上がり, 用いる }
- 算 5 = { 方法, 某, 術 }
- 算 6 = { 勘定, 数量, 数 }

“算盤 1”の1次属性に“算”がある。これで“算”は“計算 1”の1次属性と同じ表記を持つことになるので“計算 1”の“算盤”の意味は算盤 1ということにする。次に“計算 1”の中の“算定”の2次属性の意味決定だが、“算定”は多義概念ではないので意味は一意に決まる。次に“計算 1”の中の“算”であるが、“算 1”から“算 6”の中には“算”の属する1次属性、つまり“計算 1”の属性と同じ表記を持つものがないので関連度計算を使つての意味決定となる。

● 関連度での意味決定

関連度での意味決定方法は a_i の持つ意味の中で最も X に属する属性群と関連度が高い意味を決めることから、 a_i の持つ意味が3つだとしたとき X1 のすべての属性と A のすべての属性、X2 のすべての属性と A のすべての属性、X3 のすべての属性と A のすべての属性と関連度を取り、すべての値のなかで最大値を取った時の a_i の意味番号を a_i の持つ意味番号となる。

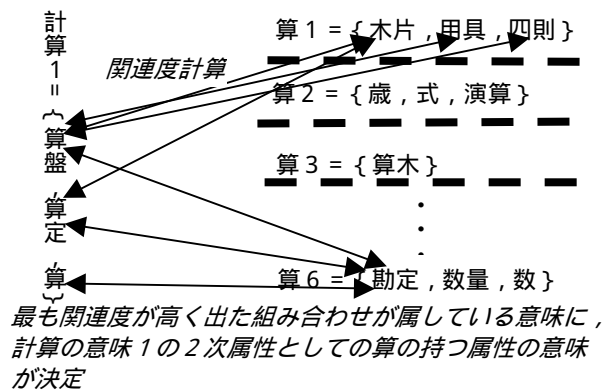


図 7 2 次属性の意味決定の関連度決定方法 (計算 1 の 1 次属性である算の意味番号を決める処理)

すべての組み合わせと関連度を取り最大値を取った“算 k”に算の2次属性がきまる。以上の処理で意味を決定する。表 4 に分割後の概念ベースの例を示す。

表 4 分割後の概念ベース (旧概念“星”の6番目の概念)

1 次属性					2 次属性
犯人	容疑者	犯罪	警察	逮捕	
犯す	被疑者	罪	国民	現行犯	
罪	犯罪	犯す	生命	自由	
現行犯	起訴	重犯	都道府県	令状	
犯罪	法律	刑罰	財産	抑留	
当人	取調	刑法	国家	警察	

4. 一義性概念ベースの評価

4.1 属性の意味分割の評価方法

分割後の概念ベースの評価方法を述べる。評価は人が判断した理想分割結果と処理の分割結果とを比較する。

評価の方法は人がサンプル概念 25 個 (全属性 2164 個) に概念構造情報を見て意味番号を振り分け処理によって決まった意味番号との一致を見た。その際、人が雑音と判断、あるいは雑音ではないが分類する意味番号が決められない属性にはそれぞれ意味番号を雑音、意味不明とした。不明と雑音と判断された属性は今回の実験では評価対象から除いた。図 8 に平均値分割と最大値分割との成功率の比較を示す。

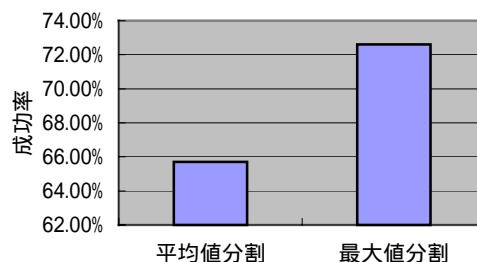


図 8 属性の意味分割の評価結果

平均値分割と最大値分割では最大値分割のほうが成功率にして 8%ほど良い分割がされている。このことから最大値分割を採用する。

4.2 評価尺度を用いた評価

● 多義性解消評価方法

多義性がどの程度解消されたかを評価するために次の評価データを尺度として作成した。ランダムに抽出した多義性を持つ概念 M_X に対し、意味 i の属性の同義語、類義語を高関連概念 M_A 、意味 j の属性の同義語、類義語を無関係概念 M_C として 63 セット作成する。表 5 に例を示す。

表 5：一義性評価尺度の例

M_X	M_A	M_C
紫	醤油	色
上	年上	位置
黒	犯人	色
丸	正解	球
顔	看板	耳
光	希望	速い
星	犯人	太陽
焼く	妬む	火

高関連概念 M_A に基準概念の持つ複数の意味から同義語を選び、無関連概念 M_C に基準概念の複数の意味の中で高関連概念とは違う意味での同義語、類義語を選んだ。これは例えば紫と醤油の関係から概念紫の意味はの場合醤油の意味であるために、色とは関係がないということ期待しての評価である。

M_X と M_A, M_C それぞれの関連度をそれぞれ r_A, r_C とする。常に $r_A \gg r_C$ となり、 r_A が 1 に近く、 r_C が 0 に近ければ、関連度計算の理想的な結果といえる。

表 6：一般評価尺度の例

M_X	M_A	M_B	M_C
樹木	木	木の葉	頭
天気	天候	雨	写真
時刻	時間	時計	消しごむ
海	海洋	波	耳
瞳	目	顔	靴
人	人間	動物	箱
子供	童	大人	雲
辞書	辞典	本	家

● 一般評価方法

連想メカニズムへの一般的な利用の面から評価するため、改良前(多義性概念ベース)と改良後(一義性概念ベース)を比較するため次のような評価データを作成する。1 セットが 4 つの概念からなる表 6 のような評価尺度を人手によって作成する。4 つの概念のうち 1 つは、そのセットの基準となる概念 M_X である。残り 3 つは、概念 M_X と同義または類義の概念 M_A 、関係のある概念 M_B 、関係のない概念 M_C である。今回使った評価尺度には、このような概念が 500 セット入っている。

M_X と M_A, M_B, M_C それぞれの関連度をそれぞれ r_A, r_B, r_C とする。 $r_A > r_B > r_C$ の結果が得られる時正解とする。

5. 結果と考察

分割することによって概念数は 33,699 個から 52,516 個に増加した。全属性数は変わらないので 1 概念が持つ平均属性数が減ったことになる。図 9 は分割後と分割前の属性数における概念数の相対度数を比較するグラフである。

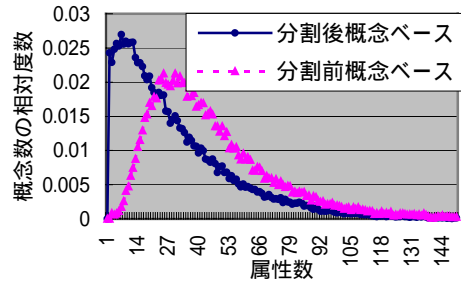


図 9 属性数の比較

分割後は属性数が少ない概念が多くなっている。

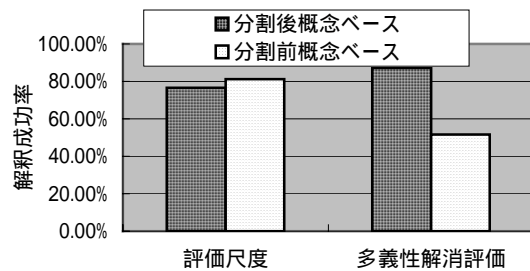


図 10 評価結果の比較

評価結果(図 10)が示すように、今回作成した分割後の一義性概念ベースは一般評価尺度においては、分割前の概念ベースに比べ 5% ほど悪い結果になった。これは一般の評価尺度では多義性を考慮にして作られていないことや、関連度計算において最も適切な属性数が 30~50 個である^[5]のに対し、分割後概念ベースでは 10 個が最も割合が多くなっていることが原因と考える。しかし、多義のために、意味分割されていない概念ベースでは大小の判定が難しい一義性評価尺度において 35% ほど良い結果が出ており、多義性が問題となる常識判断メカニズムにおいて改良されている。概念ベースを用いた連想機能の向上を図ることができる。

6. おわりに

本稿では、コンピュータに人間に近い知的な判断を行わせることを目的に、「概念ベース」の多義性の問題を解決し、一義性の「概念ベース」を構築する方法について述べた。概念ベースに必要な属性には同義や類義、他、関係のある語の多様性は必要だが、多義性は区別する必要がある。多義のために常識的な判断を狂わせる評価においてもこの概念ベースによると適切な判断ができることも示せた。

また、概念ベースには同じ意味の概念が複数あるという問題がある。これに対しては同義概念を統合する同義語統合の必要がある。その場合でも、多義のために統合できない概念に対しても今回の厳密な意味分割によって統合が可能となると考えられる。

謝辞

本研究は文部科学省からの補助を受けた同志社大学の学術フロンティアの研究の一環として行った。

参考文献

- [1] 笠原要, 松澤和光, 石川勉, 国語辞書を利用した日常語の類似性判別, 情報処理学会論文誌, Vol.38, No.7, pp.1272-1283 (1997)
- [2] 渡部広一, 河岡司, 常識的判断のための概念間の関連度評価モデル, 自然言語処理, Vol.8, No.2, pp39-54 (2001)
- [3] 井筒大志, 東村貴裕, 渡部広一, 河岡司, 概念ベースを用いた属性集合の一致度による概念間の関連度評価方式, 人工知能学会全国大会, 1A1-03 (2001)
- [4] 小島一秀, 渡部広一, 河岡司, 常識判断のための概念ベース構築法 - 国語辞書からの抽出した概念間の論理関係の利用, 同志社大学理工学研究報告, Vol.42 No.1, pp1-8, (2001)
- [5] 入江毅, 渡部広一, 河岡司, 概念ベースにおける属性数の検討と概念間の関連度計算方式, 信学技報, AI99-82, pp.37-44 (2000)