

知的メカニズムのための概念間の類似度定量化方式

1M-6

入江 肇[†] 渡部 広一[†] 河岡 司[†] 松澤 和光^{††}[†] 同志社大学工学部^{††} NTT コミュニケーション科学研究所

1. はじめに

柔軟な利用者インターフェースを持つ情報処理システムの基本は知的メカニズムであり、その中核となる機構は概念ベースと概念の類似を利用した連想機能と考えられる。本稿では、概念ベース[2]を利用し、概念間の類似度を定量的に評価する方式を提案している。本方式の特徴は類似度を概念の2次属性集合で意味として評価している点である。また、4万語の概念ベースを使用した実験により本方式が人間の感覚により適合することを示した。

2. 概念の定義

ある概念Aは、その概念の意味特徴を表す単語の集合で表し、それらを概念Aの1次属性と呼ぶ。

$$\text{概念 } A = \{a_1^1, a_2^1, \dots, a_i^1, \dots, a_N^1\}$$

また、概念Aの1次属性はさらにその1次属性(概念)の意味特徴を表す単語の集合で表せる。これらの単語集合を概念Aの2次属性と呼ぶ。

$$\text{1次属性 } a_i^1 = \{a_1^2, a_2^2, \dots, a_i^2, \dots, a_M^2\}$$

このように、概念Aはn次属性まで定義可能であるが、本稿では2次属性まで概念が定義されているものとする。

3. 概念間類似度の評価モデル

2つの概念Aと概念Bの類似度を評価する場合、観点というものが必要になってくる。例えば、「飛行機」と「自動車」の類似度と「飛行機」と「すずめ」の類似度はどちらが大きいかを考えてみる。この場合、「乗り物」という観点で見る場合と「飛ぶ」という観点で見る場合とでは人間の感覚では類似度の逆転が生じるのが普通である。したがって、概念Aと概念Bの類似度は観点と呼ぶ第3の概念Cを指定したときに計算できるとするのが自然である。

しかし一般には、観点を陽に指定せずに2つの概念間の類似度を求めたい場合も多い。このような場合は、文脈やその他の状況から観点が陰に指定されていると考えられるが、そのような陰に指定されている観点を取り出すことは、現状では困難である。

Measuring Semantic Similarity Model between Concepts for Intelligent Mechanism

Takeshi Irie[†], Hiroyasu Watabe[†], Tsukasa Kawaoka[†] and Kazumitsu Matsuzawa^{††}

[†] Faculty of Engineering, Doshisha University

^{††} NTT Communication Science Laboratories

そこで、本稿では概念Aと概念Bの類似度を

- (1) 観点が指定されていない場合
 - (2) 観点が指定されている場合
- の2つの場合に分けて取り扱う。

3.1 観点なし類似度評価アルゴリズム

概念Aと概念Bの類似度を $\text{Sim}(A, B)$ とし、アルゴリズムを以下のように定義する。

I. 概念Aの1次属性の並びを固定する。

II. 概念Bの各1次属性を対応する概念Aの各1次属性との一致度の合計が最大になるように並び替える(並び替わったBを B_x とする)。

$$B_x = \{b_{x1}^1, b_{x2}^1, \dots, b_{xi}^1, \dots, b_{xN}^1\}$$

ここで、一致度を次のように定義する。

◆一致度 $\text{Match}(a_i^m, b_j^m)$ の定義

2つのm次属性 a_i^m と b_j^m の一致度は、それぞれの $m+1$ 次属性の一致単語数を0から1の範囲に正規化したものとする。すなわち、

$$a_i^{m+1} = \{a_1^{m+1}, a_2^{m+1}, \dots, a_i^{m+1}, \dots, a_L^{m+1}\}$$

$$b_j^{m+1} = \{b_1^{m+1}, b_2^{m+1}, \dots, b_j^{m+1}, \dots, b_M^{m+1}\}$$

と表現し、 $a_i^{m+1} = b_j^{m+1}$ なる a_i^{m+1} の個数を p 個とするとき、一致度 $\text{Match}(a_i^m, b_j^m)$ を次式で定義する。

$$\text{Match}(a_i^m, b_j^m) = \frac{p(\frac{1}{L} + \frac{1}{M})}{2}$$

III. 概念Aと概念Bの1次属性数を N_A , N_B とすると概念Aと概念Bの類似度 $\text{Sim}(A, B)$ は先に定義した一致度を利用して次式のように表せる。

$$\text{Sim}(A, B) = \frac{\sum_{i=1}^{N_A} \text{Match}(a_i^1, b_{xi}^1) + \sum_{j=1}^{N_B} \text{Match}(a_i^1, b_{xi}^1)}{2}$$

ここで、1次属性同士の並び替えは文献[1]で提案している「単純法」を用いている。この方法では一致属性数の合計が最大値になるとは限らないが、比較的最適解に近い値が出来ることから、本稿で行う実験ではこの方法を採用した。なお、より最適な解を得たい場合は、遺伝的アルゴリズムなどを用いることができる[1]。

3.2 観点付き類似度評価アルゴリズム

概念Cを観点とする概念Aと概念Bの類似度を

$\text{Sim}(A, B | C)$ と書くことにする。観点概念 C で概念 A と概念 B を比較することは、観点概念 C に関する部分、すなわち観点概念 C が持つ 1 次属性に関する部分で評価するということである。アルゴリズムを以下のように定義する。

I. 観点概念 C の 1 次属性の並びを固定する。

$$C = \{c_1^1, c_2^1, \dots, c_i^1, \dots, c_N^1\}$$

II. 概念 A の各 1 次属性を対応する観点概念 C の各 1 次属性との一致度の合計が最大になるように並び替える。

$$A_x = \{a_{x1}^1, a_{x2}^1, \dots, a_{xi}^1, \dots, a_{xN}^1\}$$

III. II と同様に、概念 B の各 1 次属性を対応する観点概念 C の各 1 次属性との一致度の合計が最大になるように並び替える。

$$B_y = \{b_{y1}^1, b_{y2}^1, \dots, b_{yi}^1, \dots, b_{yN}^1\}$$

IV. 概念 A と概念 B との観点概念 C における類似度は A_x と B_y の（並びが固定された）各要素間の一致度と、観点概念 C と A_x の各要素間の一致度、および観点概念 C と B_y の各要素間の一致度の 3 者の相乗平均の（相加）平均とし、次式で表す。

$$\text{Sim}(A, B | C) =$$

$$\frac{\sum_{i=1}^N \sqrt[3]{\text{Match}(c_i^1, a_{xi}^1) \times \text{Match}(c_i^1, b_{yi}^1) \times \text{Match}(a_{xi}^1, b_{yi}^1)}}{N}$$

4. 実験結果と考察

4.1 観点なし類似度の実験

図 1 は概念：警察とその概念の同義語ないしは類似語間での類似度計算の結果であり、濃い線が提案方式で薄い線がベクトルの内積による方式である。ベクトルの内積による方式とは、文献[2], [3]など従来の類似度計算方式であり、概念を 1 次属性の重みを要素とするベクトルとみなし、2 つの概念の類似度を 2 つのベクトルのなす角の余弦で計算するものである。実験結果（図 1）から（実際の実験はさらに多くのサンプルについて行っているが、ほぼ同様な結果を得ている）、提案方式では従来の方式に比べ、10 語の類似度の差が小さくなっている様子がわかる。

ここで、同義ないしは類似の関係にある概念を対象として類似度の計算を行ったわけであるから、人間の感覚と同じようにお互いの類似度の差ができるだけ小さくなることが望ましい。そういう意味から考察すると、提案方式によって得られた結果は評価できるものと考えられる。

4.2 観点付き類似度の実験

表 1 は 3 つの概念（看護婦、婦警、医者）について、2 つの観点概念（女、病院）によって類似度を

計算した結果である。表 1 より、2 つの観点概念（女、病院）の両方について、人間の感覚と一致するような結果が得られたことがわかる。

このような実験を無作為に 37 通り行ったところ、2 通りの観点概念が両方とも人間の感覚と一致しなかった例はなかった。また 2 通りともに人間の感覚と一致したのは 19 通りで、一方のみ一致したのは 18 通りであった。よって、本実験により提案した方式は十分評価できるものと考えられる。

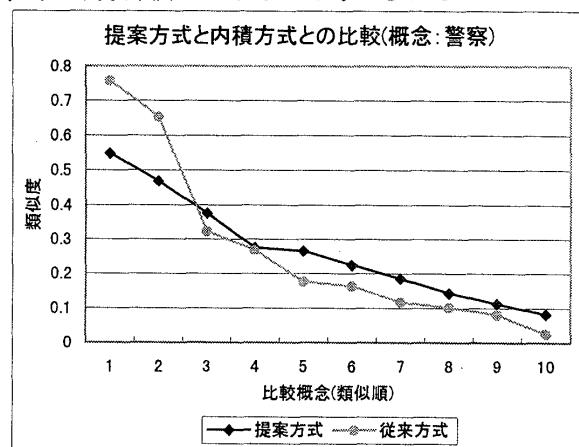


図 1 観点なし類似度実験結果

表 1 観点付き類似度実験結果

概念A	概念B	観点C	類似度
看護婦	婦警	女	0.142
看護婦	医者	女	0.038
婦警	医者	女	0.037
概念A	概念B	観点C	類似度
看護婦	婦警	病院	0.058
看護婦	医者	病院	0.222
婦警	医者	病院	0.054

5. おわりに

本稿では概念の意味を考慮し、概念の 2 次属性までを用いた概念間類似度の評価方式を提案した。今回の実験の範囲では提案方式により従来のベクトル内積方式に比べ人間の感覚により適合する結果が得られたが、今後それぞれの適用領域を明らかにする必要がある。

参考文献

- [1] 浮田知彦、渡部広一、河岡 司：概念間の関連度計算への遺伝的アルゴリズムの適用、情報処理学会春季全国大会、IU-2 (1998)
- [2] 笠原 要、松沢和光、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 (1994)
- [3] Salton, G. and McGill, M.: Introduction to modern Information Retrieval, McGraw-Hill (1983)