

Fuzzy c-means 手法を利用したコーパスからの容器状物体の形状獲得

黒澤 義明

竹澤 寿幸

広島市立大学大学院 情報科学研究科

{kurosawa, takezawa}@ls.info.hiroshima-cu.ac.jp

1. はじめに

コーパス資源の蓄積は、自然言語処理研究の対象を拡大した。これに伴い、比喩等の言語現象を扱う研究も増えている。しかし、まだ不十分であると言えよう。

例えば、「鍋を食べる」という換喩の解釈としては、文字通りの解釈はありえない。鍋は食用可能な物体ではないからである。そのため、別の解釈～『容器-中身』という関係性を媒介に、“鍋(の中身)を食べる”という解釈～を要求する。

それならば、いわゆる容器として分類可能な名詞を『容器-中身』の関係性を持つ物体と考えれば、上記の解釈は容易と考えるかもしれない。しかし、実際には、山梨(1988)も指摘しているように、「押し入れをかき回す」という例文の解釈のためにはさらなる知識～容器ではない物体も『容器-中身』の関係性を持つという知識～が必要となるからである。

すなわち、いわゆる容器としての物体でなくとも、何かを入れるという機能を持つ物体は『容器-中身』の関係性を満たす可能性を有する。ここに広い意味での「容器状物体」という新しい分類が必要となる。では、こうした容器状物体の定義を如何にコンピュータに与えるのか？ 単に閉空間を持つ物体と定義すればよいのか？

黒澤ら(2008)では、コーパスからの『容器-中身』の関係性記述を試みている。彼らは、“Aの奥”、“Aの底”等、物体Aを表す名詞とともに共起する表現に着目することを提案した。彼らの考えはシンプルであり、例えば、「鍋には深さがある一方で奥行きはさほどない。このため“鍋の底”とは言えても“鍋の奥”とは言いがたい。また、瓶は細長いため“瓶の先”と言える。しかし、鍋の場合には“鍋の先”という表現は難しい」といった人間が持つ言語直感を、コーパスから獲得し、利用することを提案したのである。

本研究も、このような共起表現(彼らが見立て詞と呼ぶ)を用いることにより、一種の言語直感のコンピュータへの構築を目指す。

ここまで述べたように、黒澤ら(2008)は上記の目的意識の上で実験を行っている。しかし、彼らの手法～Hoffman(1999)によるpLSA(probabilistic Latent Semantic Analysis)を用いた次元の縮約・整理後、Kohonen(2001)の自己組織化(Self-Organizing Map: SOM)を行う手法～には問題が残る。すなわち、最終的にSOMによるクラスタリングを行っているため、複数のクラスタへの所属が困難という点である。

そこで、本研究ではソフトクラスタリングのアルゴリズム fuzzy c-means による分類を試みる。このアルゴリズムにより、名詞が持つ複数の関係性を記述できるはずである。

2. fuzzy c-means による容器状物体の形状分類

2.1. fuzzy c-means のアルゴリズム

本研究では、Bezdek(1981)による fuzzy c-means を用いた物体の関係性記述を試みる。以下に、アルゴリズムを記す。

2.2. fuzzy c-means のアルゴリズム

n 個のデータを c 個のクラスタに分類するため、次の目的関数を考える。ここで帰属度 u_{ik} は $n \times c$ 行列、距離 d_{ik} は $d_{ik} = \|x_k - v_i\|^2$ を示す。 x_k は k 番目のデータを、 v_i は i 番目のクラスタ中心を表す。

$$J = \sum_{k=1}^n \sum_{i=1}^c (u_{ik})^m d_{ik}$$

有効なクラスタの獲得 (J の最小化) のため、以下の2式の計算(クラスタ中心及び帰属度の計算)を繰り返し行う¹。また、パラメータ m ($m > 1$) は曖昧さを決定する。 m が 1 に近づくと曖昧さは減少し、k-means 法の結果と同様になる。

¹ 導出については、宮本(1999)等、参照のこと。また、本研究は初期中心としてデータ中のランダムな1点を選んでる。

$$v_i = \sum_{k=1}^n (u_{ik})^m x_k / \sum_{k=1}^n (u_{ik})^m$$

$$u_{ik} = \left[\sum_{k=1}^n \left(\frac{d_{ik}}{d_{jk}} \right)^{1/m-1} \right]^{-1}$$

2.3. クラスタへの帰属確率

最終的に個々のデータの帰属度が確率で表現され、また、帰属度の和が 1 になるように正規化される。次に挙げる例 (表 1) では、 $c = 2$ を指定したと仮定している。

表 1 fuzzy c-means による結果例

		帰属度 (クラスタ)	
		First	Second
User	A	0.95	0.05
	B	0.65	0.35
	C	0.08	0.92
	D	0.15	0.85

このように複数のクラスタへの帰属が記述可能であるため、『容器—中身』の関連性と同時に、他の関連性についても獲得できるはずである。したがって、SOM を使用した先行研究よりも有効な結果が得られると考えられる。

3. 実験と考察

3.1. 言語データ

本研究で用いる言語データは、基本的に黒澤ら (2008, 2011) で収集されたデータと同一である。その収集・加工手続きを示す。

① 容器状名詞の収集

まず、典型的な容器だけでなく、比喩的に解釈可能な名詞 A を収集する。野口 (2005) の現代仮名遣い作品から、10 回以上登場する「A の中」という名詞句 265 個を収集した (ex. 頭の中)。なお、この手続きだけでは筒状の物体が少なくなるため、筒状の形状を持つ単語を追加した上で、さらに「光」・「音」・「流れ」等の抽象的、もしくは形状を有しない名詞を除き、152 語を実験材料とした²。

² 黒澤ら (2008) や黒澤ら (2011) とは異なる語が採用されている。

② 見立て詞毎出現率算出

池原ら (1997) の「2610 場」の下位分類から、8 語 (端, 角, 口, 奥, 先, 席, 底, 隅) を選び、容器状名詞との共起頻度を計算し、出現率にデータ変換を行った。

$$m_i \text{ の出現率} = m_i / \sum_{i=1}^n m_i$$

3.2. 実験手続き

2 章で説明した fuzzy c-means による分類を行った。なお、パラメータ m については、1.1 から 2.0 までを変化させた。また、パラメータ c についても複数の検討を行っている。しかし、紙面の都合により、 $m = 1.4$, $c = 20$ の結果についてのみ述べることにする。

3.3. 実験結果

実験により得られたクラスタの一部を示す (表 2)。ここでは、0.05 程度の値を持つ名詞のみ表示している³。なお、2 列目は便宜上の分類を表している⁴。

4. 考察

4.1. 全体の傾向について

黒澤ら (2008, 2011) 同様、基本的なクラスタリングとしては有効であると考えられる。例えば、クラスタ B においては、基本的に『文字通りの容器』が分類されている。確かに一部、容器以外の名詞 (「靴」, 「ボート」) も含まれている。しかし、その形状のみから考えると、山梨 (1995) に見られる <容器のイメージスキーマ> によって説明可能であると考えられる。

4.2. 分類の偏りについて

表 2 中の数字は、クラスタへの帰属確率を表しており、1 に近いほど、決定的な帰属性を示していると言える。したがって、1 に近い値を多数持つクラスタは、帰属の一意性が高いクラスタであると言

³ $1/c$ 程度の意。確率の分配が等価に行われた場合、各帰属度確率は $1/c$ となる。したがって、何らかの意味のある値は、この $1/c$ を超えた場合であるとえられる。

⁴ 容器や日用品の区別も明確ではない。あくまで、参考程度のカテゴリであると考えて欲しい。

える。一方、0.5以下の値を多数持つクラスタは、他のクラスタへの帰属の可能性を残したクラスタであると言える。

表2では、クラスタA及びクラスタBの一意性が高いと言える。半数以上のクラスタで0.9以上の値を有しているからである。一方、その他のクラスタでは最大値も小さく、0.5以下の値が大多数を占める。

このように、各クラスタが持つ帰属度の値は一樣ではないため、実際に使用する際（例えば、換喩の検出）には、何らかの区別が必要であることが示唆される。

4.3. ソフトクラスタ化による効果

先に述べたように、クラスタAやBでは、一意性が高いため、このようなクラスタでは、ソフトクラスタリングによる効果はあまり期待できない。その一方で、クラスタGやFについては、一定の効果が確認できる（表3）。

クラスタGもクラスタFもどちらも何らかの長さを示唆させる名詞である。ただし、前者では「松林」等、より空間的に長大な長さを示している。言い換えれば、より空間的な広がりを持っていると考えられる。その一方で、後者ではホース等の長めの物体を示しており、ある種の閉鎖性を示唆していると考えられる。

この点は、表3に記されたコーパス中の見立て詞との共起欄（特に「奥」と「先」との共起）から確認することができる。すなわち、クラスタGでは、見立て詞「奥」との共起が大きいため、奥行きが反映された空間的な広がりを感じさせる。しかし、クラスタFでは、見立て詞「奥」との共起はほとんどなく、また、見立て詞「先」との共起がクラスタFよりも大きいため、一点への焦点化が行われていると考えられる。したがって、より閉鎖的な印象を感じさせている。

このようなクラスタの空間的開放性・閉鎖性を利用することにより、例えば、比喩の生成～僕の未来はまるでパイプのようだ～に適用できると考えられる。比喩の理解の研究はあっても、比喩の生成の研究はなされていないため、こうした利用法は非常に有意義であると考えられる。

5. おわりに

本研究は、コーパスから容器状物質の形状抽出を試み、fuzzy c-meansを使用した実験を行った。今回の実験結果は、先行研究のSOMを使用した結果に比べて、複数のクラスタへの分類～特に、人間の認知的な観点による分類～を確認しており、ソフトクラスタリングアルゴリズムの有効性を確認できたと言える。

今後の課題としては、一部のクラスタに、サブクラスタとして判定できるデータも含まれていたことから、より適切な分類尺度も考慮することが挙げられる。

謝辞

この研究の一部は、平成22、23、24年度広島市立大学特定研究費（一般研究）の補助を得ている。関係各位に感謝申し上げる。

参考文献

- Hofmann, T. (1999). "Probabilistic Latent Semantic Indexing." in Proceedings of the 22nd Annual International ACM SIGIR Conference on Research and Development in Information Retrieval", pp.50-57.
- 池原悟・宮崎正弘・白井諭・横尾昭男・中岩浩巳・小倉健太郎・大山芳史・林良彦(1997) "日本語語彙大系." 岩波書店.
- Kohonen, T.(2001). "Self-Organizing Map, 3rd Edition." 徳高平蔵, 岸田悟, 藤村喜久郎訳 (2005). "自己組織化マップ." シュプリンガー・ジャパン.
- 工藤拓. "PLSI", <http://chasen.org/~taku/software/plsi/>
- 黒澤義明 (to appear in 2013). "クラスタリング手法による物体形状記述と人間の印象判定～比喩・換喩の検出に向けて～." 言語の創発と身体性, ひつじ書房.
- 黒澤義明, 竹澤寿幸 (2011). "クラスタリング手法を利用したコーパスからの容器状物体の形状獲得." 言語処理学会年次大会.
- Kurosawa, Y., Hatamoto, N., Hamada, S., and, Takezawa, T. (2012) "Comparing Clustering Algorithms for Psychomime Classification using Probabilistic Latent Semantic Analysis and Fuzzy c-Means." In Proceedings of the 16th International Conference on Knowledge-Based and Intelligent Information & Engineering Systems.
- 黒澤義明, 原章, 市村匠 (2008). "換喩検出を目的とした自己組織化マップ SOM による物体の形状マップ生成." 言葉と認知のメカニズム, pp.353-374, ひつじ書房.
- 野口英司(2005)『インターネット図書館 青空文庫』はる書房.
- 山梨正明(1988)『比喩と理解』東京大学出版会.
- 山梨正明(1995)『認知文法論』ひつじ書房.
- 山梨正明(2000)『認知言語学原理』くろしお出版

