

A Study of the Number of Analogies contained in Large Multilingual Corpora

大規模多言語コーパスにおける類推関係の推算

イヴ・ルパーシュ (Yves LEPAGE)

yves.lepage@atr.jp

国際電気通信基礎技術研究所 (ATR)

1 はじめに

言語学では、古くから「類推」が新たな文を解析や生成する方法の一つとして知られている¹。古典ギリシャやローマ時代では、「アノマリア」²という概念で表層的には類推関係にならないタイプの意味的類推関係を示した。又、逆に、意味的には類推関係が満たされないが、表層的に類推関係が成立する場合があるため (例: *Abby is baking vegan pies. : Abby is baking. :: Abby is too tasteful to pour gravy on vegan pies. : ??Abby is too tasteful to pour gravy on.*)、生成文法派の言語学者は、生成文法理論の立場から説明している。

本研究では、大規模コーパスにおける類推関係数を推定を試みた。「真の類推関係」のみを数える目的で、表層では類推関係が満たされたもののうち、同時に意味的にも類推関係になる場合を数える。表層レベルでは、提案済みの定義を採用する。意味のレベルでは、翻訳しても意味は文が異なっても伝わるという仮説を立て、ある言語での表層と意味の両方の類推関係が、他の言語へ翻訳する際にも、表層上の類推関係を満たす文として翻訳できるかを議論する。

2 コーパス

本研究では、「旅行基本表現集」(BTEC)を対象に分析を行なう。これは、旅行と観光に関する会話表現の広範囲でカバーする多言語コーパスである。音声翻訳に関する研究コンソーシアム C-STAR により³、162,318 文の表現が様々な言語に翻訳されている。本研究では、そのうち、中国語、英語と日本語のコーパスを使用する。中国語では 96,234 異なる文があるのに対し、英語では 97,769、日本語では 103,274 の文がある⁴。本コーパスの傾向を表 1 の通りに示す。BTEC の文は平均的に短いと言える。

¹主な参考文献は、(Paul 20, chap.5)、(Bloomfield 33, p.276)、(Mounin 68, p.119-120)、(Itkonen 94, p.48-50)、(Pullum 99, p.340-343)。

²ギリシア語の形容詞 *ανομάλος* から。正しい分析は *αν + ὄμαλος* (平かない、不似合い)。Lat. *similis* ↔ gr. *ὄμαλος*。そこから、英語の「anomaly」。

³<http://www.c-star.org/>。

⁴日本語の文の少しでも多い数は、漢字と平仮名での異なる書き方で説明できる (例: 「ください」又は「下さい」)。図 5 にその実際例を示す。

3 表層的類推関係

表層的類推関係を次のように形式化できる⁵。

$$A : B :: C : D \Rightarrow \begin{cases} d(A, B) = d(C, D) \\ d(A, C) = d(B, D) \\ |A|_a + |D|_a = |B|_a + |C|_a, \forall a \end{cases}$$

ここでは、 a は記号 (文章で、文字) を示し、 $A \sim D$ は記号列 (コーパスで、文) である。記号 a が記号列 A 中に出現する頻度を $|A|_a$ で表す。また、 $d(A, B)$ は A と B の間で削除を 1・挿入を 1 のコストで⁶計算した編集距離である。この式が成立する場合なり $A \sim D$ が、表層的類推関係にあると言える。

4 類推関係文の抽出

4.1 表層的類推関係

それぞれの言語 (中国語、英語、日本語) のコーパスを対象に表層的類推関係にある文数抽出した。表層的類推関係の例を図 3 で示す。また、表 1 はそれぞれの計算の結果である。これによると抽出された類推関係は大変多い。例えば、英語では約 10 万の文を含む約 250 万の類推関係が抽出できた。

図 1 の左側は、コーパスに対する類推関係数の変化を示す。

一つの文あたりいくつの類推関係が存在しているか調査した結果、英語に対し、図 2 の左グラフの通りとなった。右側は Zipf の法則に従った形になっているが、類推関係が多い文が類推関係数 750 と 2,000 の周りに不自然なピークがある。文の長さとその文に類推関係の数の検討をしても相関関係を見つかることができなかった。

4.2 意味的類推関係

自動的に意味的類推関係文を抽出するためには、意味の表示が必要なるが、検討対象コーパスには、そのような意味の表示はない。しかし、これは多言語コーパスであるので、いろいろな言語で対応している翻訳文の集合からなっている。翻訳

⁵(Lepage 01)、FG/MOL。

⁶ここは、置換や変換は数えない。

文対の中に意味の表示がかくされていると考えられる。

単言語においては、表層的類推関係のうち、どれが意味的類推関係も満たすが判定することは困難である。しかし、ある言語で表層的にも、意味的にも、類推関係にある文はそれを他の言語に翻訳した場合でも表層的類推関係を満たすという仮説をたて、その妥当性を議論する。

$$\forall \mathcal{L}, \exists A \in \mathcal{L}('A'), \\ \exists B \in \mathcal{L}('B'), \\ \exists C \in \mathcal{L}('C'), \\ \exists D \in \mathcal{L}('D'), \quad A : B :: C : D$$

$$\Rightarrow 'A' : 'B' :: 'C' : 'D'$$

すなわち、ある4つの文に対して、どの言語でも表層的類推関係を満たす翻訳が見つかり、意味的類推関係も満たす文であることが考えられる。

しかし、実際にこの仮説に基づき意味的類推関係を抽出する場合には、二つの問題がある。一つは、全ての言語でこのテストを行うことは困難である。もう一つは、ある言語でも、全ての可能な翻訳(換言)の作成は困難である。それに対して、本研究で使用するコーパスは換言文が少ない。英語の一つの文に対して中国語への可能な翻訳文は平均1.20文である。また、日本語への翻訳文は平均1.52文である。

この様な二つの問題は、英語コーパスを対象に実験を行った。その表層的類推関係、2,384,202関係の中に238,135関係は中国語でも満たされていることが分かった。その関係には25,554文が含まれている。従って、中国語訳の表層的類推関係も用いると、英語での表層的類推関係は約10%は意味的類推関係であること、即ち「真の類推関係」であると考えられる。日本語訳の表層的類推関係を用いると、共通の類推関係は336,287関係(英語での表層的類推関係の14%)であった。その関係には13,602文が含まれている。このような関係の例を図4に示す。また、中国語訳にも日本語訳にも表層的類推関係にある場合は、13,602文を含む68,164「真の類推関係」が得られた。次章でこの低い値とその理由について議論する。

図1の右グラフは抽出となる文数に対する「真の類推関係」の数の変化を表す。指数関数的な表層的類推関係の数の増加とは違って、線形的な成長となることが分かる。

5 考察

英語での表層的類推関係の数は、250万あったが、このうち、日本語に表層的類推関係が成立するものは30万関係になり、又中国語でも成立するものは7万関係になった。ここからもれた類推関係は本当に意味的類推関係ないと言えるのだろうか。実際はその中に意味的類推関係が含まれ

ていると考えられる。この原因としては、一文あたりの翻訳候補数が少ないことが考えられる。

手法を逆に考えると、即ち、二つの違う言語間で同時におこる表層的類推関係が意味的類推関係であるという仮定は類推に基づく翻訳の原理である⁷。この手法では、すくなくとも一つの表層的類推関係を持つ英文を日本語へ翻訳すると、平均174翻訳候補(換言文)が得る(正しくない文も含まれている)。実際にコーパスにある1.52文の候補と比較するのは難しい。この例を図5で示す。

また、換言文を増やすと、類推関係の可能な組み合わせが増加する。全ての対象した言語で表層的類推関係が偶然に見つかる可能性も大きくなる。それに対して、対象する言語の数を増やさなければならぬと考えられる。

6 おわりに

本論文では、実験により大規模多言語コーパスにおける類推関係の数を推定した。また、異なり10万文のコーパスで、「真の類推関係」、即ち表層でも意味でもの類推関係を推定した。その結果、約1万4千文を含む約7万の「真の類推関係」が得られることができた。本実験で用いたコーパスは、一文あたりに翻訳候補が少ないので、得られた数は非常に低い値になったと考えられる。

7 謝辞

本研究は通信・放送機構の研究委託により実施したものである。

参考文献

- Leonard BLOOMFIELD, *Language*, Holt, New York, 1933.
- Esa ITKONEN, Iconicity, analogy, and universal grammar
Journal of Pragmatics, 1994, vol. 22, pp. 37-53.
- Yves LEPAGE, Analogy and formal languages,
Proceedings of FG/MOL 2001, Helsinki, 2001.
- Yves LEPAGE, 比例類推的記号列形式言語, 言語処理学会第7回年次大会, 東京大学, 2001年3月, pp. 93-96.
- Georges MOUNIN, *Clefs pour la linguistique*, bibliothèque 10/18, Seghers, Paris, 1968.
- Hermann PAUL, *Prinzipien der Sprachgeschichte*, Niemeyer, Tübingen, 5th ed. 1920 [1st ed. 1880].
- Geoffrey K. PULLUM, Generative grammar, *The MIT Encyclopedia of Cognitive Sciences*, The MIT Press, Cambridge, 1999, p. 340-343.

⁷(Lepage 01), 言語処理学会第7回年次大会。

	類推関係の数	類推関係を 含む表現の数	コーパス中の 表現の数	文の長さ (文字数)		
				平均	±	標準偏差
中国語	3,567,123	77,039	96,234	11.00	±	5.77
英語	2,384,202	53,250	97,769	35.14	±	18.81
日本語	1,910,062	53,579	103,274	16.21	±	7.84
「真の」類推関係	68,164	13,602				

表 1: 多言語コーパスにおける類推関係の数。

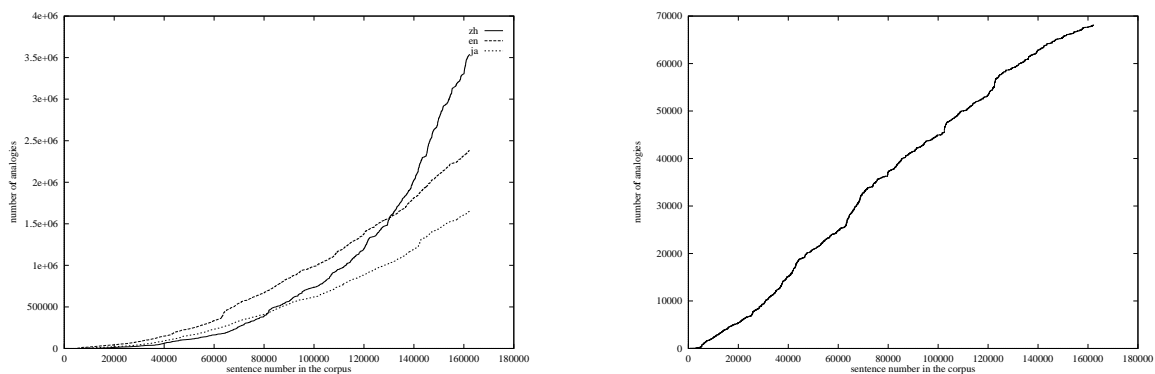


図 1: コーパス数に従い、類推関係数の変化。左：中国語と英語と日本語の表層的類推関係のみ。右：その言語で同時表層的類推関係（「真の類推関係」）数の変化。

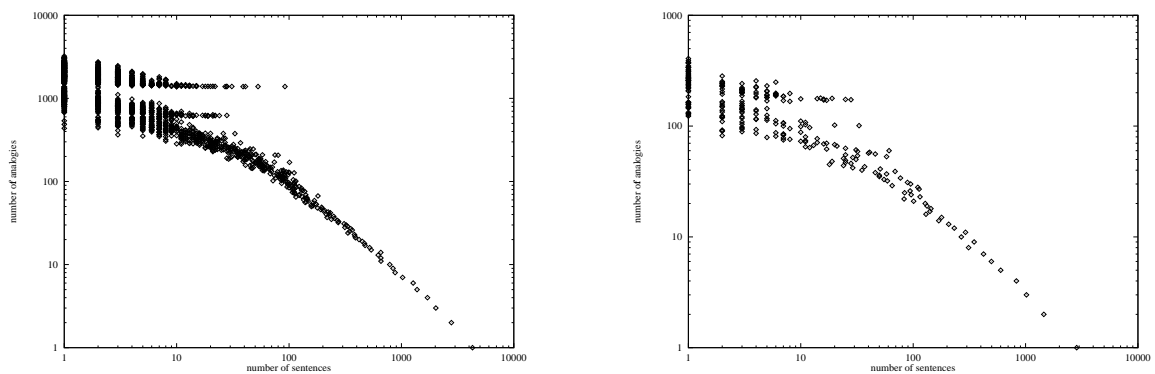


図 2: 1文あたり類推関係数の分布。左：英語の表層的類推関係のみ。右：中国語と英語と日本語の全てで層的類推関係（「真の類推関係」）数。横軸；表現の数（log スケール）。縦軸；類推関係の数（log スケール）。

美しい。	: 美味しい。	:: 何にしますか。	: 何味にしますか。
お釣りをください。	: お釣りをもらえますか。	:: 住所を教えてください。	: 住所を教えてくださいませんか。
保険をかけますか。	: 保険をかけたいのですが。	:: ヒルトンホテルに行きますか。	: ヒルトンホテルに行きたいのですが。

図 3: 日本語での表層的類推関係の例。

メキシコ料理の ほうが好きで す。	: 中華料理のほう が好きます。	:: この辺りにメキ シコ料理店はあ りますか。	: この辺りに中華 料理店はありま すか。
↓	↓	↓	↓
<i>I prefer Mexi- can food.</i>	: <i>I prefer Chi- nese food.</i>	:: <i>Is there a Mex- ican restaurant around here?</i>	: <i>Is there a Chinese rest- aurant around here?</i>

図 4: 意味的類推関係になる 2ヶ国語での表層的類推関係の例。

コーパス中	類推に基づく翻訳結果
2× 日本の新聞をください。	210× 日本語の新聞をください。
1× 日本の新聞を見せてください。	177× 日本の新聞をお願いします。
1× 日本の新聞を見せて下さい。	147× 日本語の新聞をお願いします。
	132× 日本の新聞をください。
	105× 日本語の新聞をいただきたいのですが。
	...
	3× 日本新聞はありますか。
	3× 日本新聞をいただけませんか。
	3× 日本新聞サイズをお願いします。
	3× 日本新聞何かをお願いします。
	3× 日本分の新聞をお願いします。

図 5: 左側は英文 *A Japanese newspaper, please.* のコーパスにある実際日本語翻訳。右側は同文の頻度の最高と頻度の最低の 5 文の日本語への類推に基づく翻訳結果。合わせて、異なる翻訳候補 261 文が得た。