

6F-7

類語国語辞典を介した 意味マーカ付与

田中英輝* 江原暉将** 有賀憲和*** 松田健生***

*NHK放送技術研究所 **ATR自動翻訳電話研究所 ***カテナリソース研究所

1. はじめに

筆者らは英日自動翻訳システム⁽¹⁾で動詞と名詞の共起意味選択制限による動詞訳語選択⁽²⁾に使用するための、意味マーカの設定と、翻訳辞書英語見出し語に対する日本語名詞訳語への付与方法を研究している。一般に辞書の中で、名詞の占める割合が一番大きい。多数の語義に対して、設定された意味マーカをひとつずつ付与していく作業は、非常に労力がかかる。また、意味マーカに変更を生じた時の訂正もひとつずつ行わなければならない。そこで本稿では、角川類語国語辞典⁽³⁾を仲介して、半自動的に意味マーカを日本語名詞訳語に付与する方法を提案し、その実験を行ったので報告する。

2. 意味マーカについて

意味マーカ体系は、具体物、抽象物、気象現象、補完マーカの4つを柱とした上位下位分類の形をしており全部で81個設定した。具体物の中には、ニュース分野に対応するため[CRIMINAL]{犯罪人}を設けている。また抽象名詞はその分類の難しさを軽減するために、まず、動詞の派生形を持つか、形容詞の派生形を持つかで分類した後、意味的な分類を行っている。

3. マーカ付与の方法

意味マーカの半自動的に付与は、図1に従い、類語国語辞典を仲介して、2つの対応表を作成した後行った。

英日辞書 対応1 分類番号 対応2 意味マーカ

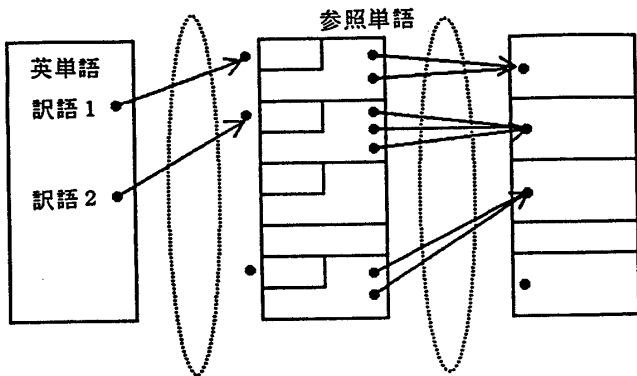


図1. 意味マーカ付与方法

3. 1. 名詞訳語と分類番号の対応

図1の対応1の部分以下の要領で作成した。角川類語国語辞典は、細分類まで含めると、2794の分類番号で単語を分類している。日本語名詞訳語に対して、その語が収録されている類語国語辞典分類番号(以下、分類番号と略す)をすべて付与した。複合語については、主要構成語を用いて番号を付与した。分類番号が見つからない語については、新たな分類番号33個を作り、それらを付与した。

3. 2. 分類番号と意味マーカの対応

図1. 対応2の部分以下の要領で、4名の作業員で作成した。

(1)すべての分類番号に対して、以下の単語を、その分類番号の参照単語として抽出した;その分類番号の項目名(整数分類番号のみ);その分類内の先頭の単語;中間付近の単語。これらの参照単語を分類番号に併記した作業ファイルを作成した。

(2)このファイルを4つに分割し、各作業員は、その内の2つについて(3)の作業を行った。

(3)各番号に付与された、参照単語(2個または3個)を見て、それに対応すると思われる意味マーカを付与する。このとき、複数の意味マーカの付与を許した。

この作業の結果、2794分類すべてに対して、2名の作業員による、意味マーカとの対応データが得られた。この2名の対応の相違を各分類について厳密に計算したところ、1600分類について相違がみられた。この相違は、抽象物の下の上位下位の間で多くみられた。

表1、表2に、同一ファイルに対する2名の意味マーカの使用頻度の違いを示した。これは、各作業員が付与したマーカ上位4つとその割合を示したものである。これらの表によれば、作業員によって、使用する意味マーカに違いがみられる。また表1、表2の作業員2に着目すると、意味マーカ使用に一定の傾向がみられる。

表1. 作業員によるマーカ付与の差1

作業員1		作業員2	
CONTACT	9.5%	ABSACT	13.5%
HUMAN	7.9%	CONTACT	13.5%
HUMSTATIC	7.6%	NOWILLCHA	7.1%
ACTION	7.2%	HUMAN	6.6%

An Assignment of Semantic Markers to Nouns through a Dictionary

Hideki TANAKA (NHK), Terumasa EHARA (ATR)
Norikazu ARUGA, Takeo MATSUDA (CLI)

表2. 作業者によるマーカー付与の差

作業者2		作業者3	
CONACT	14.1%	ACTION	17.6%
MONOACT	8.7%	HUMSTATIC	9.2%
ABSACT	8.4%	HUMAN	7.0%
HUMAN	6.4%	ADJECTIVAL	5.6%

3. 3. 名詞訳語と意味マーカーの対応

3. 1. 及び3. 2. で得られた対応表を用いると名詞訳語と意味マーカーの対応が機械的にとれる。分類番号と意味マーカーの対応は、2人分得られており、今回は、それらの和集合を対応データとした。これらの対応データを用いて、ニューステロップ英語原稿1989年9月分の中で使用頻度の高かった16個の動詞と共に名詞165語の日本語訳語に自動的に意味マーカーを付与した。その結果を検討したところ、期待された意味マーカーは、ほぼ付与されていたが、不適切な意味マーカーも付与されていた。

4. 問題

不適切な意味マーカーは大別して次の2つの理由で付与されていた。

(1) 不適切な語義の分類番号の仲介。(図2)

例えば、図3に示すようにようにbackの訳語(背)に対して、[LENGTH] (長さの概念) というマーカーが付与されていた。類語国語辞典では、{背}は104b, 123b, 601dの3カ所に登録されている。このなかで、123bは{背丈}の意味であり[LENGTH]が付与されていたが、backの語義として使用するの是不適切である。

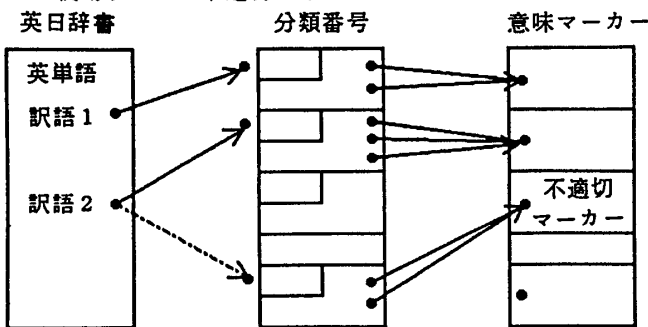


図2. 不適切語義の仲介

back

- {背後} : [RELATION]
- {背} : [RELATION]
- [LENGTH]
- [LIFEPART]

図3. backに付与された意味マーカー

(2) 類語国語辞典とマーカーとの分類不整合。(図4)

図5に示すように、advantageの訳語{強み}に対して[HUMAN] (人間) という不適切なマーカーが付与されている。類語国語辞典の中で、{強み}は、135 {強弱} 172e {長所 優れているところ}、

670c {屈強 力が優れて強いさま} の3カ所に登録されている。このうち670cの分類中に、{精锐} という単語も収録されている。分類番号と意味マーカーの対応作業の時の参照単語の一つとして、{精锐} が選ばれたため、670cに[HUMAN]も付与されてしまった。類語国語辞典の一つの分類番号が、我々の意味マーカーの複数個を含んだために発生した誤りである。

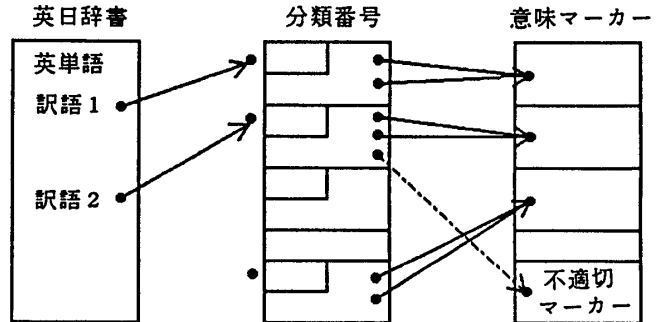


図4. 分類の不整合

advantage

- {強み} : [HUMAN] {人間},
- [WILLPROP] {形容詞派生を持つ名詞, 意志物の属性},
- [NOWILLPROP] {形容詞派生を持つ名詞, 無意志物の属性},
- [HUMSTATIC] {人間に関する抽象概念}

図5. advantageに付与された意味マーカー (一部)

5. まとめ

角川類語国語辞典を仲介して、日本語名詞訳語に、意味マーカーを半自動的に付与する手法を提案し、実験を行った結果、少ない手間で必要なマーカーを付与することができた。しかし今回の実験は、作業者が付与したすべてのマーカーを使用し、また、名詞訳語が属する分類番号すべてを利用したため、不適切なマーカーも含む結果となった。今後、この不適切な意味マーカーを削除しなくてはならない。不適切な分類番号を仲介したために付与されたマーカーは、人手で英語の訳語として適切な分類番号だけを選択するか、和英辞典を利用して、自動的に選択することで削除することが考えられる。(4) いずれにしても削除の手間は、新たにマーカーを付与する手間より少ないと思われる。また、類語国語辞典と意味マーカー分類の不整合に関する問題、作業者によるマーカー付与傾向の差の取り扱い、今後検討していきたい。

【参考文献】

- (1) 相沢, 他 「衛星放送ワールドニュースの英日機械翻訳」 情報処理学会第40回全国大会
- (2) 中瀬 「英日機械翻訳システムにおける解析手法について」 情処研資 NL69-7, 1988
- (3) 大野, 浜西 「類語国語辞典」 角川書店
- (4) 田中, 他 「翻訳辞書からの中間概念の自動抽出に関する基礎的考察」 情処研資 NL72-3, 1989