

図解辞書と LDOCE の分野コードに基づく 場面知識による英語名詞の多義性解消

角 田 達 彦[†] 羽 柴 正 輝^{††}

本稿では、図解辞書と LDOCE 電子化版の分野コードを組み合わせることによって場面に関する知識を構築し、文中の場面に関連する英語名詞の多義性解消を行う方法を提案する。場面は段落ごとにあらかじめ決定されているとし、その段落内の名詞全体に場面知識を適用する。場面に関する知識は、名詞の語義を列挙したものと、それに基づいて場面に特有な分野を集計したものから成る。場面に関連する名詞に対しては積極的に多義性解消を行い、場面に関連しない名詞に対しては、推定の誤りを少なくする必要がある。提案手法を物語文の中の名詞に適用した結果、場面に関連する名詞に対し約 88 % の再現率および約 78 % の適合率という結果が、そして場面に関連しない名詞に対し約 90 % の再現率という結果が得られた。

Disambiguation of Noun Sense by Scene Knowledge Based on Pictorial Dictionary and LDOCE Subject-codes

TATSUHIKO TSUNODA[†] and MASATERU HASHIBA^{††}

We propose a method of English noun disambiguation using scene information. We hypothesize here that each scene of each paragraph is determined beforehand. According to the scene, our method applies scene knowledge to all the nouns in the paragraph. The scene knowledge is constructed in advance using the nouns from the pictures in OXFORD-DUDEN Pictorial English Dictionary (OPED). The knowledge has two types of representation: (a) pairs of nouns and their senses, and (b) list of subject-codes of Longman Dictionary Of Contemporary English (LDOCE) accompanied by the meanings in (a). We applied our method to nouns in a narrative story; for nouns which were related to the scene the recall ratio was about 88 %, and the precision was about 78 %; for nouns which were not related to the scene the recall ratio was about 90 %.

1. はじめに

語義の曖昧性の解消は、解決困難な問題の 1 つである¹⁾。一般には、Katz らの選択制限²⁾などのように、格フレームのスロットに意味的な制限を設け、制約によって曖昧性を解消することが行われることが多い。そのためには各動詞に対する格のとりうる条件を整備することが必要となるが、それは大変手間のかかる作業である。さらに実際の文章では、読み手の記憶を反映した簡潔な表現にするために、文脈を積極的に利用することが多く、格のとりうる条件も文脈や状況によって変化することも考慮に入れる必要がある。しかし現在のところ、そのような文脈や状況の記述方法が

不明確であるし、格フレームの意味的制限をさまざまな状況に応じて獲得するという手法は見い出されておらず、人間が個別に記述せざるをえないという問題がある³⁾。

そこで、このような構造的な格関係を利用することは辞書や知識の整備を待つことにし、文章中の語と語の共起関係など現段階で比較的容易に扱える情報によって曖昧性を解消する研究が近年行われている。

その 1 つとして、コーパス中にある語と語が共起する頻度をとらえ、一種の文脈とする方法があげられる。たとえば Yarowsky らは、語の各語義と共起する語を集め、精度良く多義性を解消する手法を提案している⁴⁾。その方法ではまず、それぞれの語義に対して共起しやすいと思われる語を人間がいくつか最初に与え、その語と元の調べたい語がともに現れる文章を大量に集める。その中で元の語義とより共起性の高い語を取り出し、精度を向上させる。この方法は、コーパ

[†] 京都大学工学部
Faculty of Engineering, Kyoto University

^{††} 日立製作所
Hitachi, Ltd.

スから文脈的知識を獲得する方法を示したものであり、精度も良く、意義がある。しかし彼の研究では“plant (1. 植物, 2. 会社)”など、語義を2つと仮定し、比較的多義性解消が容易と思われる単語12例を扱っているが、実際の語では互いに区別のしにくい語義が多く存在する語も多い。そのような区別の難しい語義ごとに共起語を人間が最初に手で与える必要があり、その労力は無視できないと思われるし、見つけ出して与えることが可能かどうかは今のところ不明である。また、かなり大量のコーパスを用意し、すべての語の1つ1つの語義に対して学習させていく必要もある。

もう1つの方向として、既存の人間用の辞書を知識源として積極的に用い、最初から高いカバー率を狙うことが考えられる^{1),3)}。人間が使うために書かれた辞書は語の意味や使われ方を端的に示した1つのコーパスであり、最初からかなり良質の知識が集められたものであると位置づけることもできる。そのような辞書の使い方の1つは、語義の定義文に現れる語を用い、各語義と文中の語の共起関係をとらえる方法である。たとえば Lesk は、与えられた文章中の語を英英辞典で引き、定義文に現れる語をすべての語のすべての語義に対して集め、語の重なるの多い語義を求めることによって多義性を解消するという方法を提案している⁵⁾。この方法の有効性の詳細は、本稿の実験と考察の章で比較検討するが、概していえば、簡単な方法でありながら比較的精度の良い結果が得られる。だが、定義文はすでにさまざまな知識を持つ人間に説明しようとするためやや偏った語が現れることがあり、かつ簡潔に説明しようとするために記述量が少なく、定義文の間での重なりが見られず誤ってしまう傾向がある。また実際の文章では、その場に現れる物体などの具体的な意味を持つ語だけでなく抽象的な意味を持つ語も多く現れる。さらに複雑な例としては、“thing”など、具体的な物体を示すのか抽象的な「物事」という意味を示すかを判定する必要がある語もある。Lesk らの方法ではこのような抽象的な意味を持つ語に対して誤った意味を推定してしまうことが多いことが問題となる。

また、Guthrie らは語義の分野を用いる方法を提案している⁶⁾。英英辞典 Longman Dictionary Of Contemporary English (LDOCE)⁷⁾の電子化版につけられている分野コードを手がかりに、それぞれの分野ごとに語の定義文を集め、その中に現れる語の共起頻度をあらかじめ計算しておく。多義性解消のときは、文中の対象となる語の各語義の分野を調べ、その分野の中で、その対象となる語と共起しやすい語を上の特徴をもとに列挙する。そしてそれらが入力文のまわり

の単語と重なる数を数え、重なるの多いものを出力する。重なりが少なすぎる場合には、分野内での共起語のさらに共起語を調べ、共起語の範囲を大きくしていく。複数の候補が出てきた場合には、各語義の定義文と文章中の語の重なるの多い方を選択する手法も併用する。Guthrie らの論文では定量的な評価はなされていないが、本稿での実験での比較検討の結果は芳しくなく、上の Lesk らの方法よりも概して悪くなった。その大きな原因は、同じ分野であっても、定義文間の特徴的な語の共起はかなり少ないためである。したがって“have”や“make”, “place”など、さまざまな状況で使われる語が相対的に多くなってしまふ。実際の文章中では、このような一般的な語が数多く存在するためにノイズとなってしまふという問題と、文脈として手がかりとなる特徴的な語の種類が多く、Guthrie らの方法ではとらえきれないという問題がある。

これに対して本研究では、抽象的な語の意味を推定することは、慣用句や格関係の詳細な記述や知識の整備が進むことによって解決可能であると考え、より具体的な意味を持つ語を先に積極的に取り出し、多義性解消を行う方法を提案する。具体的な意味を持つ語は、文中の場所格などの手がかりによって特定しやすい⁸⁾場面という単位をもとに特定する。多義性解消の対象とする語は英語の名詞とする。典型的な場面に関する知識を近似するため、図解辞書の図版に現れる名詞の語義や、それが使われる分野を図版ごとに集める。図解辞書は、OXFORD-DUDEN Pictorial English Dictionary (OPED)⁹⁾を用い、語義は上記のLDOCEによって定義する。分野の割当てには、LDOCEの電子化版の、各語義につけられた分野コードを用いる。このような場面知識をあらかじめ構築しておき、その知識により、文章中に現れる名詞の語義を推定する。

ここでは、文中の場面が段落ごとに同定できたと仮定し、その段落の中の全名詞に知識を適用する。このとき、場面に関連する名詞の多義性が解消でき、場面に関連しない名詞に対する誤りが少ないことが望ましい。図解辞書とシソーラスを用いて語義検索の優先づけを行う手法¹⁰⁾がすでに提案されているが、ここでは語義と分野づけのために、電子的に利用可能であるLDOCEを用い、また場面に関連しない名詞に対しても実験および評価を行ったことが独自の点である。提案手法を物語文に適用し、これらの評価を行う。

2. 場面知識の構築

場面に関する知識は、図解辞書 OPED に基づいて作る。この辞書は、物の名前を、それが現れる場面や、

kitchen	
1.	housewife
2.	refrigerator (fridge, <i>Am.</i> icebox)
3.	refrigerator shelf
4.	salad drawer
:	:
44.	kitchen table

図1 図版の物体に割り当てられた名前のリストの例

Fig. 1 A list of names assigned to objects in a picture.

物の形などの種類から人間が引くことを目的に出版されたものであり、日常生活に出てくる場面が384枚の図版に網羅的に描かれている。それぞれの場面に出てくる典型的な物体には、それを表す語が割り当てられており、リストにまとめられた形で図版につけられている。たとえば台所(kitchen)の図版の場合には、図1のように、'housewife'から'kitchen table'まで44個の物体の名前があげられており、それぞれの番号が図版中の対応する物体につけられている。

本研究では、この図版全体を人間が場面ととらえているものの近似であるとし、登場する物を、場面に現れやすい物であると仮定する。登場する物の名前は、名詞によって表現されているものがほとんどであるため、これらの名詞のみを知識源として扱うことにする。ここで、図1にある'kitchen table'を見ると分かるように、図版中の名詞には修飾語と被修飾語があるが、両方ともその場面に密接に関係があると判断したため、区別をせず用いることにした。

知識の種類としては、これらの名詞のそれぞれに語義を割り当てたものと、さらに図版に現れる語の語義の意味的なまとまり(分野)をとらえたものと、2つを構築することにする。後者の意味的なまとまりは、LDOCEの電子化版で各語義につけられている「分野コード」によってとらえる。分野はFU(家具:Furniture)など、100種類ほどあり、特定の分野の中で使われると判断された語義に対して割り当てられたものである。複数の分野に関連すると判断されたものに対して複数のコードが割り当てられている語義もあれば、特に特定の分野に関連すると判断されず何も割り当てられていない語義もある。本研究でこのような分野の違いによる意味的なまとまりを知識にする理由は、図版にはたまたま現れなかったが場面に関連すると思われる語(未登録語)に対しても、多義性解消の処理が行えるようにするためである。

場面知識の構築は、図2のような手順で行う。たとえば、図解辞書の「台所」の図版には、机が描かれ、それには'table'という語がつけられている。この場

- | | |
|----|--|
| A. | 図解辞書の各図版に現れる名詞を列挙する。 |
| B. | 上記Aの各名詞に語義を割り付ける。
LDOCE中の語義の中から、図版中で使われていると思われる語義を、人手で選ぶ。 |
| C. | 上記Bで得られた各語義についての分野コード(LDOCE電子化版中)を調べる。
図版ごとに、現れる分野コードをまとめる。 |

図2 場面知識構築の手順

Fig. 2 The method of constructing the scene knowledge.

[台所]	語義:	[table, 1. a piece of furniture...], (1)	
		[tea, 3. a hot brown drink...],	
		...	
	分野:	[HH, FO, FU, EG, BV, CE,	(2)
		HR, BO, PM]]	
[寝室]	語義:	[drawer, 2. a container...],	
		[bed, 1. a piece of furniture...],	
		...	
	分野:	[HH, FU, CL, TE, BE, HR,	
		AF, EG, SI]]	

図3 場面知識の2つの表現: 語義フレームと分野フレーム

Fig. 3 Two types of representing the scene knowledge: Sense frame and subject frame.

合の'table'の語義を人間が判断すると、LDOCE中では、'1. a piece of furniture...'であるため、これらを対にして「台所」の場面の「語義フレーム」に蓄える(手順B)。また、この語義1をLDOCE電子化版で調べると、'1. [FU] a piece of furniture...'のように、分野コードFUが割り当てられているので、この分野コードを「台所」の場面の「分野フレーム」に加える(手順C)。以上を図版中のすべての名詞に対して行い、「語義フレーム」と「分野フレーム」のそれぞれをまとめる。その結果、図3のような知識が得られた。それぞれの分野フレーム中のコードの表すものは、HHはhousehold(家にある物)、FOはfood(食べ物)、FUはfurniture(家具)、HRはhour(時間)、EGはengineering(工学)、BOはbotany(植物)、CLはclothing(衣類)、TEはtextiles(織物)、BEはbeauty(化粧品)、AFはart(美術品)SIはscience(科学)であり⁷⁾、分野は頻度順に並べてある。この頻度を用いて多義性解消の細かい制御を行うことも考えられるが、今回は評価が複雑になるため、行わないことにした。これらの他の場面に対する知識も同様に、それぞれ別々に作る。

3. 場面情報による多義性解消

多義性解消は、図4のような手順に従って行われる。まず、E. Brillの開発したPart-Of-Speech Taggerと、

1. 形態素解析により名詞を取り出し, 原形化.
2. 場面知識の語義フレーム ([名詞, 語義] の対) の中に, 調査中の名詞が見つければ, それに対応する語義を出力し, 次の名詞を調べる.
↓ 見つからないとき
3. LDOCE で各語義の分野コードを調べ,
(a) 場面知識の分野フレーム ([分野 1,...]) に登録されている分野を持つ語義を, すべて出力する.
(b) 分野フレームに登録されている分野を持つ語義が1つもなければ, すべての語義を出力する.

図4 入力文中の名詞の多義性解消の手順
Fig. 4 The method of noun disambiguation.

Penn 大学で開発された morph を用いて入力文の形態素解析を行い, 名詞を取り出す. それらの名詞の語義を, 上で構築した場面知識によって推定する. ここでは, LDOCE 中の語義の中から選ぶことを, 語義の決定と定義する.

場面知識の適用であるが, まず, 場面知識の語義フレームの中で, 処理対象の名詞と同じ名詞を探す. 見つければ, その語義を出力する. たとえば, 文章中の 'table' という名詞の多義性を解消したいとする. その段落は台所の場面であったとする. 2章「場面知識の構築」で作った, 台所の知識の語義フレーム (1) を見ると, [table, 1. a piece of furniture...] というスロットがあるため, '1. a piece of furniture...' をそのまま出力する. もしも, この語義フレームに対象としている名詞が登録されておらず, 見つからなかった場合は, 次に分野フレームを探す. たとえば, 'preserve' という名詞は図解辞書になかったため, 語義フレームには登録されていない. しかし, 台所の知識の分野フレーム (2) に FO という分野 (食べ物: food) があり, かつ, 'preserve' の語義の中に FO という分野コードが割り当てられている語義 (砂糖菓子) があるので, その語義を出力する. 各単語で, 場面知識の持つ分野と一致する分野がつけられた語義が複数ありうるが, それらの語義をすべて出力する. これに対し, 'thing' などの抽象名詞のように, 分野フレームにも照合する分野を持つ語義が1つもない場合は, ここではその名詞は処理の対象としないと考え, そのすべての語義を出力し, 他の処理にゆだねる.

4. 物語文を対象とした実験と評価

4.1 評価対象とした文章の中の名詞の性質

今回提案した多義性解消手法を評価するための文章は, モンゴメリー作『赤毛のアン』の英語原作より, 台所の場面を含む7段落と, 寝室の場面を含む7段落

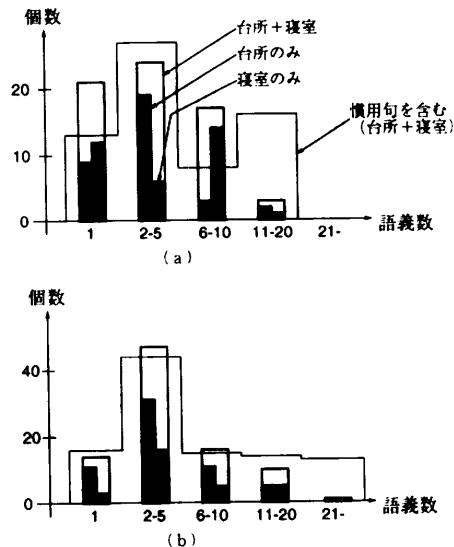


図5 評価の対象とした名詞の語義数の分布: (a) 場面に
関連する名詞の場合, (b) 場面に
関連しない名詞の場合
Fig. 5 Distribution of the number of senses, (a) of nouns
related to each scene, and (b) of nouns non-related
to each scene.

を人間が判断して切り出したものである.

その文章を形態素解析した結果, 台所の文章の中より 109 語, および寝室の文章の中より 84 語の, あわせて 193 語が名詞と判定された. そのうち, 人間が判断して名詞であり, かつ LDOCE 中に正解と判断できる語義が存在したものは, 台所の場面では 98 語, 寝室の場面では 68 語の, あわせて 166 語であった. この中でさらに, 'one' や 'three' など品詞の解釈すら安定していない語を除き, 台所の場面に現れる 90 語と寝室の場面に現れる 63 語の, あわせて 153 語を実験および評価に用いた.

これら評価対象となる名詞の語義数を調べた結果, 図5のように, 語義が1つのものは約2割あり, 25個の語義 (慣用句を入れると49個) を持つものまで, 幅広く分布している. 全体の語義数の平均は4.6 (慣用句を入れると7.5) 個であった. ただし, 慣用句を同定することは別の手法で扱うものとし, 今回の研究の対象外とした. これらの名詞を人間が判断し, それぞれの場面に関連する語と, 関連しない語とに分けて評価を行う. 場面に関連する語の典型的な例は, 「台所」では 'dish (皿)' であり, 「寝室」では 'bed (ベッド)' である. 一方, これらの場面に関連しない語の典型的な例は, 'thing (物事)' や 'face (顔)' などである. 今回対象とした名詞153個のうち, 場面に関連する語は65個 (42.5%) で, 場面に関連しない語は88個 (57.5%) であった.

図5から分かるように, 場面に関連しない語の方が,

場面に関連する語よりも、語義数が多いという傾向がある。これは、場面に関連しない語の方が抽象的な名詞が多く、使われ方が多様であるために、語義数が多くなっているためである。しかし、場面に関連する語にも、語義数の多いものが多くある。このため、場面に関連する語の多義性を解消することは意味がある。

4.2 評価方法とその目的

多義性解消の場合、人間が正解と判断する語義は、各名詞に対して複数ありうる。また、多義性解消処理からも、複数の解が出力される。このため、各名詞に対して次のような再現率と適合率を定義する。

$$\text{再現率} = \frac{\text{処理の出力のうち正解の数}}{\text{人間が正解と判断する語義の数}} \quad (1)$$

$$\text{適合率} = \frac{\text{処理の出力のうち正解の数}}{\text{処理が出力する語義の数}} \quad (2)$$

さらにここでは、複数の名詞をまとめて評価するため、式(1)、(2)のそれぞれの平均を求める。以下、再現率と適合率は、それぞれ平均をさすものとする。

本研究の目的は、場面に関連する名詞の多義性解消である。したがって、場面に関連する名詞に対する処理の結果は、再現率と適合率がともに高い方がよい。一方、場面に関連しない名詞に対しては、ここでの処理によって生じる誤りがなるべく少ない方がよい。再現率が高い方がよい。だが適合率は、低くともよい。場面に関連しない名詞に対しては、無理に語義を絞り込むことをせず、全語義を残しておき、他の処理で多義性を解消することを想定しているためである。

4.3 比較を行った手法について

本稿で提案する手法と比較するため、次の3つの手法も同時に実験を行った。

- (1) 各名詞の語義として LDOCE 中で先頭に登録されている語義を出力する
- (2) Lesk の手法⁵⁾
- (3) Guthrie らの手法⁶⁾

このうち、(1)の手法も検討する理由は、LDOCE は、語義が日常使われやすい順に登録されているからである。したがって、他に知識がない場合には、先頭の語義を解と見なす方法は妥当である。(2)の手法は、前述のように、文章中の各語の各語義の定義文を列挙し、互いに語の重なりが多い定義文を持つ語義を選択する方法である。定義文中のどのような種類の語を数えればよいかは任意性があるが、Lesk は定義文中の自立語のみの重なりを数えている。今回の語で実験した結果ではさらに絞り、名詞と動詞のみを列挙し重なりを数えるにせし、さらに、'be'、'make'、'use' などの、頻繁に使われノイズとなりうる動詞を外して処理を行っ

表 1 評価した名詞全体に対する多義性解消の結果

Table 1 Evaluation of all nouns.

	全体 166 個	LDOCE の先頭語義	Lesk の手法	Guthrie らの手法	今回提案 する手法
再 現 率	台所 90 個 寝室 63 個 平均	68.7 % 75.4 % 71.5 %	49.8 % 56.4 % 52.5 %	53.3 % 54.8 % 53.9 %	86.7 % 92.9 % 89.3 %
適 合 率	台所 90 個 寝室 63 個 平均	70.4 % 79.4 % 74.1 %	50.6 % 57.2 % 53.3 %	45.4 % 41.4 % 43.8 %	56.8 % 61.0 % 58.5 %

た結果が最も良かったため、これを比較対象にすることにした。(3)の手法は、LDOCE の分野コードごとにその分野コードの割り振られた語義を集め、それぞれの定義文の中で語の共起を調べておき、その共起関係を対象とする語と文章中の他の語との間に適用する。語の分野ごとの重なりは、Guthrie らの研究と同じ 2 語とし、それに満たない場合は、近接語を適用する処理を繰り返す。このようにして文章中の語と最も共起性の高い語義を求めて出力する。そして同じ分野で複数の語義が候補として残った場合は、それらの定義文と、対象とする文章との語の重なりを調べ、多い方を出力する。この場合も、'be'、'make'、'use' などの語がノイズとなったため、外すことにした。手法 2、3 とも、重なりを数える文章の大きさは段落とした。

4.4 多義性解消の結果と考察

(A) 対象名詞全体に対する多義性解消の結果

まず、対象とした名詞全部に対して多義性解消を行った結果を示す。評価した名詞は 153 個である。そのうち、台所の場面を含む段落に現れた名詞は 90 個、寝室の場面を含む段落に現れた名詞は 63 個であった。これらの名詞を評価した結果、表 1 に示されるように、提案手法では、全体の再現率は 89.3%であった。一方、先頭の語義をとる手法では、全体の再現率は 71.5%、Lesk らの手法では 52.5%、そして Guthrie らの手法では 53.9%であった。この結果から、提案手法では、少なくとも誤りを出力する割合が小さいことが分かる。適合率に関しては、提案手法では 58.5%、先頭の語義をとる手法では 74.1%、Lesk らの手法では 53.3%、Guthrie らの手法では 43.8%という結果が得られた。

提案手法の適合率が低くなっているのは、以降の評価で分かるように、場面に関連しない名詞に対して、無理な絞り込みをせず、すべての語義を残しておくという、安全な結果になったからである。したがって、この数値が低いことは、提案手法が本質的に悪いことを示すものではない。また、先頭の語義をとる手法の

表2 場面に関連する名詞に対する多義性解消の結果
Table 2 Evaluation of nouns related to each scene.

場面関連 65 個	LDOCE の先頭語義	Lesk の手法	Guthrie らの手法	今回提案 する手法
再 現 率	台所 32 個 77.6 %	71.4 %	59.4 %	81.3 %
	寝室 33 個 84.8 %	65.2 %	53.0 %	94.0 %
	平均 81.3 %	68.3 %	56.2 %	87.7 %
適 合 率	台所 32 個 81.3 %	70.3 %	51.3 %	72.7 %
	寝室 33 個 87.9 %	66.7 %	51.0 %	82.6 %
	平均 84.7 %	68.5 %	51.1 %	77.7 %

再現率と適合率をともに悪くする原因となっているものは、場面に関連しない名詞に関しても語義を1つに絞り込み、その結果誤ったためである。

以降、これらの名詞を人間が判断し、場面に関連する名詞と関連しない名詞とに分け、詳細な評価を行う。手法の比較に関しては、先頭の語義をとる手法との比較を主に行うことにし、LeskやGuthrieの手法の比較と考察は次章で論じる。

(B) 場面に関連する名詞に対する多義性解消の結果

上記の(A)の153個のうち、場面に関連すると人間が判断した名詞を取り出し、評価した結果を示す。場面に関連する名詞は65個である。そのうち、台所の場面を含む段落に現れた名詞は32個、寝室の場面を含む段落に現れた名詞は33個であった。これらの名詞を評価した結果、表2に示されるように、提案手法の再現率の平均は87.7%であり、適合率の平均は77.7%であった。一方、先頭の語義をとる手法の再現率の平均は81.3%であり、適合率の平均は84.7%であった。これらは本研究で扱う場面の情報が有効であるかを直接評価するものであるため、大変重要な指標となる。

提案手法で、図解辞書に載っていて、語義を与えたために多義性解消ができたものは、台所での'door', 'table', 'tea', 'dish', 'dinner', そして寝室での'table', 'bed', 'pillow', 'chair'であった。これらは辞書中での意味のとおり文章中で使われたものである。

次に、図解辞書にはないが、分野コードが合ったために多義性解消ができたものは、台所での'preserve', 'cake', 'couch', また寝室での'frill', 'pin', 'bed'であった。たとえば、'preserve'という名詞には次の3つの意味がある。

- (1) [FO] (often in comb.) a substance made from fruit boiled in sugar, used esp. for spreading on bread; JAM
- (2) [HF] a stretch of land or water kept for private hunting or fishing

- (3) something considered to belong to or be for the use of only a certain person or limited number of people (EX) She considers the arranging of flowers in the church to be her own preserve

[]内は分野を示し、FOは「食べ物」、HFは「狩猟・魚釣り」を表す。図3のように、台所での分野フレームにFOが入っているため、上の中から砂糖菓子を示す語義(1)が正しく選ばれた。

図解辞書にはないが、分野コードが合ったために、部分的に多義性を絞れたものは、台所ではtray(3つから2つへ)、寝室ではmat(4つから3つへ)であった。それらの分野はどちらも家の用品を示すHHであった。除外することのできた語義は、それぞれ、「デスクトレイ」と「(髪などの)もつれ」であった。

提案手法で誤った例は、台所では'dish', 'chamber', 'dinner', 'meal', 寝室では'wall', 'light'であった。この中で、'dish'と'dinner'は、上記のうまくいった例にも含まれている。これは、同じ語が似たような状況で異なる意味で使われているためであり、本手法の境界点の1つを示している。たとえば'dish'は図解辞書では「皿」という意味で載っており、ここでもそのまま適用したが、実際の文章の文脈では「料理」という意味が正しいため、誤った。また'dinner'は辞書では「夕食の食べ物」という意味であったが、実際には「食事という状況」という意味で使われていた。このような語義は、本来の意味から派生的に使われだしたものが多。

また、台所での'bread'と'butter', 寝室での'floor', 'corner'などは、場面に関連すると思われるのに積極的に絞り込めなかった。'bread'には「パン」, 「食べ物」, 「生活の糧」, 「お金」の4つの意味が、そして'butter'には「バター」と「ペースト」の2つの意味、'floor'には8つの意味、'corner'には6つの意味がそれぞれあるが、どれも分野コードが割り当てられていないため、多義性解消ができなかった。この場合には、他の語義を誤って支持することもないため再現率は高いが、多義性解消できる語がその分少なくなるため、適合率が下がる。LDOCEの先頭の語義をとる手法に比べて提案手法の適合率が低いのは、このためである。

先頭の語義をとる手法であるが、寝室の場面では、かなり高い再現率と適合率を示している。これは、寝室にある物は、あまり多義なものがなく、先頭の語義が示す典型的な意味で用いられていることが多いからである。一方、台所の場面では、「食べる」という行為

表3 場面に関連しない名詞に対する多義性解消の結果
Table 3 Evaluation of nouns not related to each scene.

場面非関連 101 個	LDOCE の先頭語義	Lesk の手法	Guthrie らの手法	今回提案 する手法
再 現 率	台所 58 個 63.8 % 寝室 30 個 65.0 % 平均 64.2 %	37.9 % 46.7 % 40.9 %	50.0 % 56.7 % 52.3 %	89.7 % 91.7 % 90.4 %
適 合 率	台所 58 個 67.2 % 寝室 30 個 70.0 % 平均 68.2 %	39.7 % 46.7 % 42.1 %	42.1 % 30.9 % 38.3 %	48.0 % 37.2 % 44.3 %

が行われるために、文章中に出てくる語には、本来の意味から派生した意味を持つために誤ったものがあった。それらは台所では、'tea' (先頭語義では「お茶の葉」であるが、文章中では「紅茶」の意味で使われた)、『dish』(「皿」→「料理」)、『dinner』(「夕食の食べ物」→「夕食」という状況)、『meal』(「食べ物」→「食事」)であった。寝室では、『wall』(「防御壁」→「部屋の壁」)、『corner』(「へり」→「隅」)、『light』(「ランプ」→「ロウソク」)であった。

(C) 場面に関連しない名詞に対する多義性解消の結果

上記の(A)の153個のうち、場面に関連しないと人間が判断した名詞を取り出し、評価した結果を示す。場面に関連しない名詞は88個である。そのうち、台所の場面を含む段落に現れた名詞は58個、寝室の場面を含む段落に現れた名詞は30個であった。これらの名詞を評価した結果、表3に示されるように、提案手法では、再現率は90.4%であった。一方、先頭の語義をとる手法では、全体の再現率は64.2%であった。この結果から、提案手法では、少なくとも誤りを出力する割合が小さいことが分かる。適合率に関しては、提案手法では44.3%、先頭の語義をとる手法では68.2%という結果が得られているが、この値は重要でない。提案手法の適合率が低くなっているのは、無理な絞り込みをせず、すべての語義を残しておくという、安全な結果になったからである。したがって、この数値が低いことは、提案手法が本質的に悪いことを示すものではない。場面に関連しない名詞は、扱う対象とせず、無理な絞り込みを行わずに、そのまま素通りさせたいという、本研究の目的に合っている。

提案手法で、図解辞書に載っていて、語義を与えたために多義性解消ができたものは、台所での'table'だけであった。これは同じ段落内で台所の場面から廊下の場面に移動し、そこに'table'があったため、偶然うまくいった例である。

次に、図解辞書にはないが、分野コードが合ったために多義性解消ができたものは、台所での'birch'、

'mare'、'room'、また寝室での'minute'であった。たとえば、ここでの'birch'は木の一種でPMという分野コードを持つ。台所の図版にも'plant'という語があり、その分野コードがPMであるために選択することができた。

図解辞書にはないが、分野コードが合ったために、部分的に多義性を絞れたという例はなかった。

提案手法で誤ったものは、台所では、『parlour』(場面知識:「飲物」→文章中:「部屋」)、『collar』(「機械のパーツのリング」→「シャツの襟」)、『glass』(「砂時計」)、『グラスコップ』、『眼鏡』→「ガラスの材質(の)☆」)などであり、人間でも判断に迷うものが多い。その原因は、BV(飲物)、EG(工業製品)、HH(家庭用品)という分野のものが台所の図版にあったためである(特にEGを持つものは「電線(lead)」であった)。寝室では、『point』(場面知識:「温度計の目盛」→文章中:「針などの先」)、『eye』(「ホックの穴」→「人間の目」)などを誤った。原因は、SI(科学)、CL(衣類)という分野のものが寝室の図版にあったためである(特にSIを持つものは「鏡(mirror)」であった)。

先頭の語義をとる手法では、台所では'appearance'、'east'、'yard'、'flood'、'hollow'、'thing'、'world'、'home'、'company'(その他多数)、そして寝室では'consolation'、'shiver'、'indication'(その他多数)などの、もともと多義・抽象的・派生的であるものの語義を推定しようとしたため、誤ったものがほとんどである。本研究で提案する手法では、これらの語を推定することを自動的に避けることができたため、再現率が高いという望ましい結果が得られた。

5. 手法の全体の考察

5.1 比較した各手法の考察

まずLDOCEの最初の語義を出すという方法であるが、単純でありながらかなり良い結果が得られた。再現率は本稿で提案している手法ほど高くないものの、適合率は7割程度の良い値を安定して保っている。よく使われる語義を出してみるという方法は、処理を1段階で終わらせ、結果を出すだけの場合は、ある程度の精度が得られることを示している。しかし、再現率が落ちるため、複数の処理で少しずつ曖昧性を解消することに用いる場合には絞り込み過ぎの傾向がある。ただし、尤度を計算するための情報の1つとして使う

☆ これは複合語を作り修飾することが仮定されている語義を持つ名詞である。

場合には、大変有効に働く可能性が大きい。

次に Lesk の手法であるが、場面に関連した語に関しては、本稿で提案している方法や LDOCE の最初の語義を出力する方法に比べると良くはないが、ある程度の精度が得られている。それに対し場面に関連しない語に関して多くの語の語義の推定に失敗している。この傾向は、文章中で意味のまとまりを作りやすい場面のような状況をとらえてはいるが、抽象的な関係をとらえることには失敗しており、誤った語義を出力してしまっていることを示している。また、誤りの大きな原因は、定義文が必ずしも実際の意味のまとまりをとらえる手がかりとして有効でない場合があることである。人間が読んで区別をするために書かれたものであるため、語義そのものではなく、他の語義との違いが書かれているものも少なくない。また yard の語義の1つに「裏庭」があるが、単に 'BACKYARD' という言い換えが書かれている場合もある。このようなまで対処するには、かなり複雑な機構が必要となる。

最後に Guthrie らの方法は、われわれの予想に反してかなり悪い結果になった。その最大の原因は、定義文に使われる語は、実際の文章のさまざまな語とは照合しないということである。特に LDOCE は基本単語のみで定義文を記述しており、語彙が限られている問題がある。これは他の一般の辞書にも当てはまる傾向と思われる。この傾向が強いため、分野ごとに定義文中で共起する単語には、特徴的な語を取り出すための計算をしているにもかかわらず、'put' や 'get' などのさまざまな状況で使われる語が多く混じり込み、判別を困難にする。たとえば、家具を示す FU という分野で 'table' と共起する単語のリストは、上から順に leg, make, support, low, movable, be, sit, cupboard, serve, pairs, ... となるが、実際の文章では、これらの語はほとんど出てこない[☆]。これに対し、食べ物を示す FO という分野で 'table' と共起する単語のリストは、get, stand, somewhere, napkin, guest, base, sit, price, menu, accord, ... であり、get や guest, sit などが文章中の語と照合してしまうことにより、'table' の意味を誤って推定してしまう^{☆☆}。さらに彼らの方法では、それでも語義を絞り込めない場合には、定義文と文章中の語の直接の共起を求めているが、一般的に使われる語を数えてしまい誤る率がこれによってさらに高まっている。

Lesk の手法も、Guthrie らの手法も、一般の英英辞

典の語義の定義文を利用しているために、語の集め方の難しさが生じている。1つの解決策は、意味そのものをより機械向けに書いた辞書やシソーラスを利用することである。本研究では、そのようなものを良質に近似する1つの知識源として、図解辞書を用いた。

5.2 本手法の限界点と拡張性

本研究では、視覚的な場面に関連する名詞の多義性を解消し、場面に関連しない名詞は無理に多義性解消せず次の処理にゆだねることを目的にした。実際にこの手法を用いる場合、以下のような点を考慮する必要がある。

場面には今回のような視覚的な場面だけではなく、湯気のような絵として描きにくいものの場面や、季節によって変化する場面、社会的な場面、また動作主の動作に関わる場面など、多くの場面が考えられる。今回の手法はこのような視覚的でない場面を扱うことはできず、限界点となる。しかし、そのような比較的抽象的な場面や状況も、今回のような典型的な視覚的場面に結び付いて行われることが多い。たとえば先生が生徒に質問し生徒が黒板に答えを書くといった一連の状況は、学校の教室という場所との結び付きが大変強いと思われる。そのため、このような状況的な場面を知識化するには、コーパス中で文章が指す典型的な場面を特定し、その場面の中で状況に応じた行為や物体の知識を取り出すことが考えられる。視覚的な場面を表す知識は、そのための種として使える可能性が大きい。これには、知識の量を増やす側面と、異なる種類の知識を作る側面がある。質的な拡張を行う場合には、その場面にあるすべての語を集めるのではなく、行為のみ取り出すなど、新たな観点を入れる必要があると思われる。またメタファなどを積極的に検出していくことも重要であろう。

次に、本手法ではこれらの視覚的でない場面に関する語とともに、より抽象的な語に対しても積極的に多義性解消をすることをしていない。これらの語の多義性解消は大変難しいが、いくつかの手法によって解決することが考えられる。1つは、典型的な言い回しを整備していくことである。実際に今回用いた辞書 LDOCE にも慣用句がある程度記述されており、意味を特定するのに有効である^{☆☆☆}。たとえば 'note' という語はノートや記憶など、7種類の通常の意味と5つの慣用句があるが、'mental note' という組合せでは1つの慣用句「思い出」に特定できる。このような慣用句は辞書にあげられている以外にも数多くあり、そ

[☆] be, make などはノイズとなるので、扱わない。

^{☆☆} この場合の意味は、複合語をなして「テーブルの上にある～」となる。

^{☆☆☆} 本研究では対象外とした。

れを整備する 1 つの方法として N-gram のような手法¹¹⁾を用いてコーパス中から統計的に獲得することが考えられる。上述の Yarowsky らの手法⁴⁾は、共起する語をより精密にしつつコーパスから学習するという点で、同じ方向を目指しているといえる。また 'thing (物, 物事, 衣類,...)' の多義性解消のように、どのような状況が解消に役立つかも明確にされていないものもあり、統計的手法が有効であるか、解析と適用の試行錯誤による推論方法や知識の整備が有効であるかを明確にする必要がある。

5.3 場面の同定について

本稿で提案した手法は、段落ごとに場面が同定されていることを仮定したものである。場面の同定は、最初の場面に入る（あるいは切り替わる）部分の同定と、その場面が継続するか否かの判定によってなされる。現在のところ、場面切替は前置詞と名詞など表層的な手がかりにより 7 割程度が検出できるという知見が得られており⁸⁾、場面による語彙的結束性¹²⁾がより精緻にとらえられるようになれば、この精度がさらに向上するものと思われる。この詳細は別稿に譲る。

また現在のところ必要な場面の数は不確定であるが、今回対象とした物語文の原本全体で 50 程度（分けられるもの約 20、連続しているもの約 30）であり、知識源として用いた辞書 OPED の図版ですべてカバーできた⁸⁾ため、384 という図版の数は、かなり妥当であると思われる。これ以外の場面の扱いについては、これらの図版を複数組み合わせることによって対処する方法や、上記のように視覚的場面以外の場面に質的に拡張していくことなどが考えられる。

6. おわりに

図解辞書と LDOCE 電子化版の分野コードを組み合わせることによって場面知識を構築し、場面に関連する英語名詞の多義性解消を行う方法を提案した。場面知識は、名詞の語義を列挙したものと、それに基づいて場面に特有な分野を集計したものから成る。物語文の中の名詞に適用した結果、場面に関連する名詞に対し再現率約 88% および適合率約 78% という結果が、そして場面に関連しない名詞に対し再現率約 90% という結果が得られた。以上から、場面に関連する名詞の多義性解消を行い、場面に関連しない名詞の語義の推定の誤りを少なくするという目的に、本研究で提案した手法が適うものであることが確かめられた。文中の場面の同定と、視覚的場面以外の場面の扱い、そして場面に関連しない名詞の多義性解消が今後の課題である。

参考文献

- 1) 奥村 学：自然言語処理の意味的曖昧性の解消法，人工知能学会誌，Vol.10, No.3, pp.332-339 (1995).
- 2) Katz, J. and Fodor, J.: The Structure of a Semantic Theory, *Language*, Vol.39, No.2, pp.170-210 (1963).
- 3) 長尾 眞, 佐藤理史, 黒橋禎夫, 角田達彦：自然言語処理，岩波講座ソフトウェア科学第 15 巻，岩波書店 (1996).
- 4) Yarowsky, D.: Unsupervised Word Sense Disambiguation Rivaling Supervised Methods, *34th Annual Meeting of the Association for Computational Linguistics*, pp.189-196 (1995).
- 5) Lesk, M.: Automatic Sense Disambiguation Using Machine Readable Dictionaries: How to Tell a Pine Cone from an Ice Cream Cone, *ACM SIGDOC Conference*, pp.24-26 (1986).
- 6) Guthrie, J., Guthrie, L., Wilks, Y. and Aidinejad, H.: Subject-dependent Co-occurrence and Word Sense Disambiguation, *29th Annual Meeting of the Association for Computational Linguistics*, pp.146-152 (1991).
- 7) Longman Dictionary of Contemporary English, Longman Group (1978).
- 8) 角田達彦, 田中英彦：談話解析に基づく場面の同定とその評価，第 9 回人工知能学会全国大会，pp.503-506 (1995).
- 9) THE OXFORD-DUDEN Pictorial English Dictionary, Oxford University Press/日本出版貿易 (1981).
- 10) 角田達彦, 田中英彦：英語名詞の多義性解消における文脈としての場面情報の評価，自然言語処理，Vol.3, No.1, pp.3-27 (1996).
- 11) Nagao, M. and Mori, S.: A New Method of N-gram Statistics for Large Number of n and Automatic Extraction of Words and Phrases from Large Text Data of Japanese, *Proc. COLING-94*, pp.611-615 (1994).
- 12) Tsunoda, T. and Tanaka H.: Analysis of Scene Identification Ability of Associative Memory with Pictorial Dictionary, *Proc. COLING-94* (1994).

(平成 8 年 4 月 25 日受付)

(平成 8 年 11 月 7 日採録)

**角田 達彦 (正会員)**

1967年生。1989年東京大学理学部物理学科卒業。1995年同大学工学系大学院博士課程修了。工学博士。同年京都大学工学研究科助手，現在に至る。IJCNN '93 Student Award

受賞。1994年情報処理学会学術奨励賞受賞。岩波ソフトウェア科学第15巻「自然言語処理」(岩波書店，共著)。言語処理学会，情報処理学会，日本神経回路学会，電子情報通信学会，人工知能学会，日本認知科学会各会員。

**羽柴 正輝 (正会員)**

1971年生。1994年京都大学工学部電気工学第二工学科卒業。1996年同大学大学院修士課程修了。同年日立製作所入社現在に至る。情報処理学会，言語処理学会各会員。