

Title	文脈ネットワークを用いた語の多義性解消
Sub Title	Word sense disambiguation using a contextual semantic network
Author	岡本, 潤(Okamoto, Jun) 石崎, 俊(Ishizaki, Shun)
Publisher	慶應義塾大学湘南藤沢学会
Publication year	2012
Jtitle	Keio SFC journal Vol.12, No.1 (2012.) ,p.97- 111
JaLC DOI	10.14991/003.00120001-0097
Abstract	自然言語処理の高度な課題の一つである語の多義性解消のために、連想概念辞書を用いる文脈ネットワークモデルを提案する。これは概念の連想関係と概念間の定量的な距離情報を用いて入力文から文脈ネットワークを作成し、活性拡散により語義を決定する。従来のWordNetを用いる方法と比較し、本提案手法が高い正解率を得た。次の動的な文脈ネットワークモデルは、人間が文を読むように文頭から入力語を理解していく仕組みで、多義が解消できるまで順番に文脈ネットワークを拡大させていく多義性解消モデルであり、有効性を確認している。
Notes	自由論題 研究論文
Genre	Journal Article
URL	https://koara.lib.keio.ac.jp/xoonips/modules/xoonips/detail.php?koara_id=0402-1201-0007

慶應義塾大学学術情報リポジトリ(KOARA)に掲載されているコンテンツの著作権は、それぞれの著作者、学会または出版社/発行者に帰属し、その権利は著作権法によって保護されています。引用にあたっては、著作権法を遵守してご利用ください。

The copyrights of content available on the Keio Associated Repository of Academic resources (KOARA) belong to the respective authors, academic societies, or publishers/issuers, and these rights are protected by the Japanese Copyright Act. When quoting the content, please follow the Japanese copyright act.

文脈ネットワークを用いた 語の多義性解消

Word Sense Disambiguation Using a Contextual Semantic Network

岡本 潤

嘉悦大学ビジネス創造学部専任講師 / 慶應義塾大学 SFC 研究所上席所員 (訪問)

Jun Okamoto

Assistant Professor, Department of Business Innovation, Kaetsu University/
Senior Visiting Researcher, Keio Research Institute at SFC

石崎 俊

慶應義塾大学環境情報学部教授

Shun Ishizaki

Professor, Faculty of Environment and Information Studies, Keio University

自然言語処理の高度な課題の一つである語の多義性解消のために、連想概念辞書を用いる文脈ネットワークモデルを提案する。これは概念の連想関係と概念間の定量的な距離情報を用いて入力文から文脈ネットワークを作成し、活性拡散により語義を決定する。従来の WordNet を用いる方法と比較し、本提案手法が高い正解率を得た。次の動的な文脈ネットワークモデルは、人間が文を読むように文頭から入力語を理解していく仕組みで、多義が解消できるまで順番に文脈ネットワークを拡大させていく多義性解消モデルであり、有効性を確認している。

We proposed a method of word sense disambiguation using Contextual Semantic Network (CSN). This network is constructed using concept dictionaries (Associative Concept Dictionary and WordNet) which include semantic relations among concepts. In the network, an interactive activation method is used to identify meanings of a homographic ideogram on the CSN where the activation values on the network are calculated using the distance information between concepts. Dynamic Contextual Network Model is applied to the network, where the network structure changes according to the successive input words in the sentences. The model enhances the network dynamically by adding words from the input sentence until the model obtains a certain threshold to decide the appropriate meaning of the homographic ideograms.

Keywords: 連想概念辞書、多義性解消、文脈理解、動的な文脈ネットワークモデル、活性拡散

1 はじめに

自然言語をコンピュータで処理し自動要約や多義性の解消などの高度な言語処理を実現するためには、言語学的情報に基づいて構文解析や表層的意味

解析を行うだけでなく、人間が言語理解に用いている一般的な知識や当該分野の背景知識などの必要な知識(記憶)を整理し、自然言語処理技術として利用可能な形にモデル化することが重要である。

一般性のある自然言語理解のために、現実世界で成り立つ知識を構造化した知識ベースが必要であり、そのためには人間がどのように言葉を理解しているかを調べる必要があると考えている。そこで我々は、大規模な連想実験を実施し、人間の語に関する連想をまとめ構造化した連想概念辞書を構築した。(岡本 & 石崎, 2001)

また文章をコンピュータで解析をするとき、語の意味を正しく把握することが必要である。たとえば、音声合成分野での同表記異義語の「読み」および「語義」の曖昧性の解消においては、例えば「額」は文脈により『金額』や『額縁』の「がく」や『おでこ』の「ひたい」などが考えられる。そこで本研究では文脈を利用した語の多義性を解消するために、文脈ネットワークモデルを提案する。このモデルは、まず多義語を含む文全体に関する文脈ネットワークを作成する。ネットワークはニューラルネットワークを模して作成され、連想概念辞書を背景的な常識として使用する。次にネットワーク内で、ニューロンの活性値を使用して語の多義性を解消する。また、人間が文を読み進める時の言語理解過程を扱う動的な文脈ネットワークモデルを提案する。このモデルでも連想概念辞書を利用し、文から順番に入力する語に応じてニューラルネットワークを更新しながら、語に対応するニューロンの活性値を使用して語義を決定する。

1.1 人間の記憶モデルと概念辞書

初期の知識に関する研究では、人間の記憶モデルの1つとして意味的に関係のある概念をリンクで結んだ意味ネットワークモデルが提案されている。Collins と Loftus は、階層的ネットワークモデル (Collins & Quillian, 1969) を改良し、意味的距離の考えを取り入れ活性拡散モデル (Collins & Loftus, 1975) を提案した。意味的距離をリンクの長さで表し、概念間で意味的類似性が高いものは短いリンクで結び、低いものは長いリンクで結んでいる。このモデルによって文の真偽判定に関する心理実験や典型性理論 (Rosch & Mervis, 1975) について説明した。

自然言語をコンピュータで扱う時の大規模な知識

ベースの例として、電子化辞書があげられる。日本ではコンピュータ用電子化辞書として EDR 電子化辞書 (EDR, 1990) が構築されている。WordNet (Miller, et al., 1993) は George A. Miller が中心となって構築した電子化シソーラスで、人間の記憶に基づいて心理学的見地から構造化されている。EDR 電子化辞書や WordNet は自然言語処理分野などでもよく用いられている。

1.2 語の多義性解消と既存手法

辞書を利用した多義性の解消には、Collins English Dictionary における見出し語とその語義文を利用したものがある (Veronis & Ide, 1990)。これは辞書からネットワークを作成し文脈に応じた語義の選択を行っている。具体的には語義文の「word」とその語義の ID を示す「sense」をノードとして大規模なネットワークを作る。語が多義語で、複数の意味がある場合はその語義「sense」同士は抑制リンクで結ぶ。このネットワークモデルを利用し、2つの語を入力語として与え、各ノードの活性値を計算することで活性値の高いノードの語義を採用する方法を取っている。たとえば、多義語「pen」は「筆記用具のペン」、「動物を入れる檻」、「雌の白鳥」などの意味があるが、「pen と goat」が与えられた場合、「mammal」や「animal」などが発火し、「an enclosure in which domestic animals are kept.」の語義が適切であると判断するものである。

また、多義語を含む文の理解において、ネットワーク表現を用いた Massively Parallel Parsing (Waltz & Pollack, 1985) などがある。これは、入力文について統語的な情報や概念情報などをノードとしてネットワークを作り、活性拡散を用いてノードの値を計算することで特定の意味や解釈を表すノードのみが高い活性値となり、多義文の理解をするものである。日本語の多義性解消に関する研究では、シンプルベイズ法、決定リスト法、サポートベクターマシンなど機械学習手法を融合し高い精度を得ており、解析に必要な素性についてのその有効性や特徴について述べている (村田他, 2003)。

2 連想概念辞書

2.1 主な特徴

我々は(岡本 & 石崎, 2001)において、小学生が学習する基本語彙の中で名詞を刺激語として連想実験を行い、人間が日常利用している知識を連想概念辞書として構造化した。また「刺激語」と「被験者に呈示した刺激語から連想された語(連想語)」の2つの概念間の距離の定量化を行なった。従来の概念辞書は木構造で表現され、概念のつながりは明示されているが距離は定量化されておらず、概念間の枝の数を合計するなどのような木構造の粒度に依存したアドホックなものであった。人間の記憶に関する研究や自然言語処理や情報検索などに応用する際に、概念間の距離を定量化したデータベースが有用になってくると考えている。従来の概念辞書には概念に関する上位および下位概念や、部分-全体、概念の特徴や類義語などの記述、また「動作」に関して格情報を記したものはあるが、概念同士が密接に関連する「環境」を記述しているものは少ない(岡本 & 石崎, 2001)。また連想概念辞書は、表層に表れない文章中の語に関する知識(常識)などを用いた文書要約(岡本 & 石崎, 2003)(Okamoto, et. al., 2011)や比喩理解システム(坂口, 2010)に応用されている。この比喩理解システムは、連想概念辞書を利用しニューラルネットワークを用いて「課長は鬼だ」という比喩表現を「課長はとても怖い」と理解して、出力するものである。

2.2 連想実験

被験者に呈示する刺激語として、光村図書出版株式会社の「語彙指導の方法」(甲斐, 1995)に記載されている小学校の学習基本語彙の名詞から「果物」、「野菜」、「木」、「乗り物」、「家具」、「人間」などを中心として3~4階層をなす上位および下位概念の語や、「語彙指導の方法」において頻出頻度が高い基本語彙(名詞)、連想実験で連想頻度が高い語(名詞)、メタファー理解、多義語の解消など数々の研究に用いる名詞を選択した。現在の刺激語数は約1100語である。連想実験は自由連想ではなく、被験者に名詞を刺激語として呈示して、「上位概念」、

「下位概念」、「部分・材料概念」、「属性概念」、「類義概念」、「動作概念」、「環境概念」の7課題に関して連想させ、任意個の連想語をキーボード入力させる。被験者数は1刺激語に対し50人である。

また、本論文で利用する連想概念辞書は、連想実験によって得られたデータを元に概念間の距離を定量化している。刺激語 x と連想語 y との概念間の距離 $D(x, y)$ は、連想実験から得られる連想頻度 $F(x, y)$ と連想順位 $S(x, y)$ のパラメータによる線形結合で表現し、線形計画法を用いて(1)式のように最適解が求められている(岡本 & 石崎, 2001)。ここではパラメータをもとに境界条件を距離 $D(x, y)$ の値が最大で10.0程度、最小で1.0程度になるように定め、シンプレックス法によって係数の最適解を計算した。

$$D(x, y) = 0.81F(x, y) + 0.27S(x, y), \quad (1)$$

$$F(x, y) = Nx / (n + \delta), \quad \delta = Nx / 10 - 1, \quad (Nx \geq 10),$$

$$S(x, y) = 1/n \sum_{i=1}^{n_{xy}} S_{xyi},$$

$F(x, y)$ は刺激語 x から連想語 y を連想した被験者の割合、 $S(x, y)$ は連想語 y が連想された順位を平均した値、 N_x は刺激語 x を呈示した被験者数、 n_{xy} は刺激語 x のとき連想語 y を連想した人数($n \geq 1$)、 S_{xyi} は被験者 i が刺激語 x の時に連想語 y を連想した順位である。(1)式では連想頻度 $F(x, y)$ の係数が連想順位 $S(x, y)$ の係数より大きく、連想人数が概念間の距離に与える影響は大きくなっている。

図1は刺激語「紙」についての連想概念辞書の記述例である。たとえば「紙」の上位概念として「文房具」が連想されており、「0.18」は『頻度』(連想者数を被験者数で割った値)、「1.00」は『連想順位』、「3.39」は「紙」と「文房具」の『概念間距離』である。「下位概念」、「部分材料概念」、「属性概念」、「類義概念」、「動作概念」、「環境概念」なども同じ形式で記述してある。多くの被験者が同一の語を連想している場合は連想頻度が高くなり、その連想語は刺激語にとって連想しやすい語であると考えられ、概念間の距離も短くなる。

			[頻度](連想者数を被験者数で割った値)	[連想順位]	[概念間距離]
紙	上位概念	文房具	0.18	1.00	3.39
紙	上位概念	物	0.40	2.00	7.29
紙	下位概念	折り紙	0.22	2.00	3.24
紙	下位概念	トイレトペーパー	0.26	3.46	3.32
紙	部分材料概念	木	0.62	1.19	1.48
紙	部分材料概念	パルプ	0.44	3.00	2.37
紙	部分材料概念	繊維	0.40	4.15	2.81
紙	属性概念	薄い	0.62	1.36	1.52
紙	属性概念	白い	0.50	2.16	1.98
紙	類義概念	ペーパー	0.52	1.04	4.37
紙	動作概念	書く	0.64	1.66	1.57
紙	動作概念	折る	0.44	2.55	2.25
紙	環境概念	学校	0.30	1.93	2.65
紙	環境概念	机	0.28	1.57	2.67

図1 刺激語「紙」についての連想概念辞書の記述例（連想語は一部のみ表示）

2.3 既存概念辞書との比較 —規模とネットワーク密度—

2.3.1 連想概念辞書と WordNet の規模の比較

連想概念辞書では、連想された語を別の連想実験で刺激語として被験者に呈示する場合がある。WordNet は多重継承などもありラティス構造になっている。また、WordNet の概念 (synset) の ID に日本語の記述を対応させた日本語 WordNet (Isahara, et al., 2008) が作られている。連想概念辞書の刺激語数は 1096 語で、WordNet は名詞の見出し語数は約 12000 語である。また日本語 WordNet は WordNet の概念 (synset) の ID に対応する日本語を表記の揺れなども含めて割り当てているので見出し語数が多くなり、名詞のみでも 66 万語になる。連想概念辞書、WordNet、日本語 WordNet の名詞に関しての規模は表 1 のようになる。連想概念辞書の概念数は名詞のみに絞ると、連想語全体と比較して 16000 語ほど少なくなる。また連想概念辞書の概念数に対するリンク数の割合と WordNet の概念数

に対するリンク数の割合を比較すると、連想概念辞書の方が WordNet よりも多くのリンクが張られていることが分かる。

2.3.2 連想概念辞書と WordNet のネットワーク密度の比較

ここでは、連想概念辞書や WordNet 全体での概念 (synset) のつながりを一つのネットワークとみなし、そのネットワーク密度について比較を行う。また、連想概念辞書中の刺激語のみを対象としたネットワーク密度と刺激語に対応する日本語 WordNet の synset のみを対象としてネットワーク密度を計算し比較を行う。連想概念辞書の刺激語に対応する synset のみを抽出して作成した WordNet のネットワークは、刺激語に対応する synset から最上位概念である entity まで順に上位概念を辿ったネットワークに、各ノード (synset) から関係子を 1 つ分辿った synset を加えて作成する。

ネットワーク密度は、ネットワークに含まれる

表 1 連想概念辞書と WordNet の規模の比較

連想概念辞書(連想語全体)		連想概念辞書(名詞のみ)		WordNet(名詞のみ)		日本語 WordNet(名詞のみ)	
概念数	62729	概念数	46400	synset 数	82115	synset 数	51736
リンク数	274822	リンク数	176833	リンク数	231535	リンク数	163262

表 2 連想概念辞書と WordNet のネットワーク密度の比較

連想概念辞書(全体の概念)	WordNet
0.82×10^{-4}	0.34×10^{-4}
連想概念辞書(刺激語のみを対象)	WordNet(刺激語に対応する synsetのみを対象)
0.04	0.66×10^{-4}

リンクの数によって決まる。ノード同士を結ぶリンクが多いと密度が高くなり、リンクが少ないと密度は低くなる。概念のネットワークを有向グラフとしてとらえ、ネットワーク密度 D を以下の式で計算する。

$$D = \frac{m}{n(n-1)}, \quad (2)$$

ここで、 m はネットワーク中のリンクの数、 n はノード数である。

表 2 は、連想概念辞書と WordNet でのネットワーク密度の比較と、連想概念辞書中の刺激語のみを対象とした場合と刺激語に対応する日本語 WordNet の synset のみを対象とした場合のネットワーク密度を示す。WordNet はラティス構造になっているとはいえ木構造に近いのでネットワーク密度は低いと考えられる。連想概念辞書では、下位概念として様々な語を連想する場合が多くその連想頻度は低い。つまり、複数の刺激語から同一の語を下位概念として連想することは少ないので、連想概念辞書内のすべての概念についてネットワーク密度を算出すると値は低くなると考えられる。一方、連想概念辞書は刺激語から連想される語も刺激語として連想実験をしているため、刺激語に含まれる語のみでネットワークを構築すると多くの語がお互いにつながっており、ネットワーク密度を計算すると、WordNet よりも約 1000 倍ほど密なネットワークとなっていると考えられる。

3 文脈ネットワークモデルによる多義性の解消

本研究では、同表記異音異義語を対象に語の多義性解消を行う。たとえば「額」は、「かく」と読む場合は絵を飾る額縁という意味を持ち、また金額を表す意味となる場合もある。さらに「ひたい」と読む場合には、「前額部」いわば「おでこ」の意味を表す。このように日本語には同表記で違う読み・意味の語が数多くあり、たとえば「金(きん)」「金(かね)」や「札(ふだ)」「札(さつ)」などがある。

本章では、文中の語に関する文脈の意味ネットワークである Contextual Semantic Network の概要と、ニューラルネットワークでの活性値の計算方法と語の多義性の解消について述べる。また、連想概念辞書と WordNet の両概念辞書を利用して Contextual Semantic Network を作成し、多義性の解消についての精度を比較する。

3.1 文脈ネットワーク

人間は、多義語の適切な意味を前後の文脈から把握していると考えられる。本研究では多義語を含む文を入力文として、文中の語に関する文脈の意味ネットワーク (Contextual Semantic Network) を作成し、文中の表層の語のみだけでなく、その語に対する背景知識・常識・どのような状況において用いられる語であるかなどの情報をネットワークに含めることによって、その文脈に合った適切な語義を判断する。

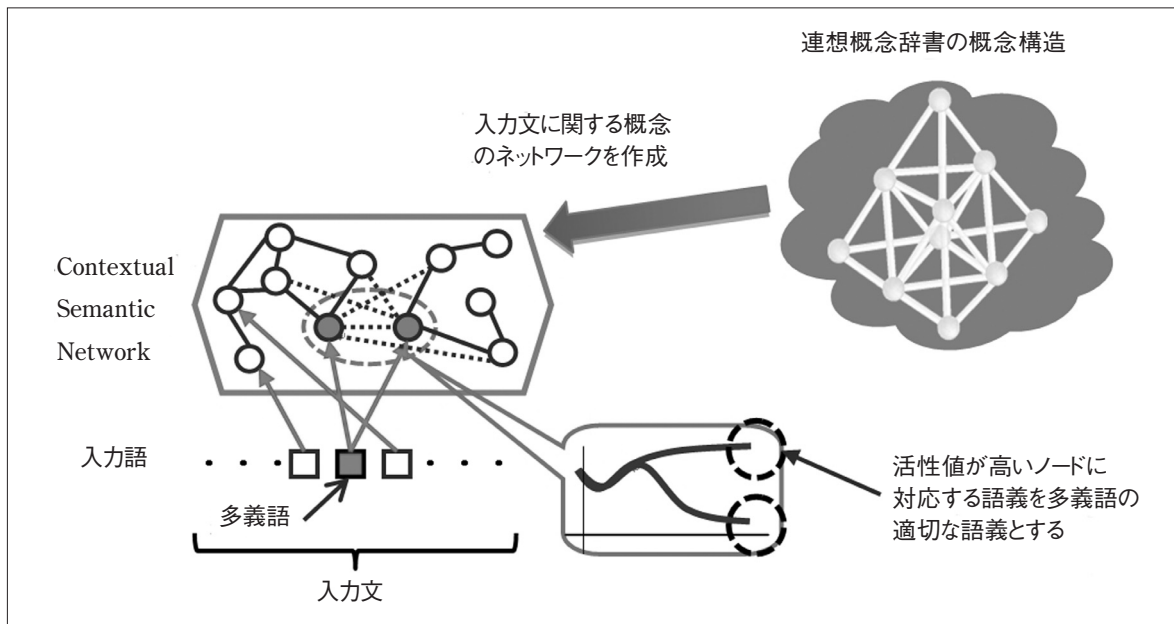


図2 Contextual Semantic Network による多義性の解消

図2は、連想概念辞書をもとに作成する Contextual Semantic Network を用いた多義性の解消の概念図である。概念間の距離が定量化された連想概念辞書のデータを用い、刺激語と連想語とその2つの概念間の連想関係と連想距離から概念のネットワークを構築する。入力文中の自立語をもとに、概念のネットワークから入力語に関する部分を中心に概念関係を抽出し Contextual Semantic Network を作成する。多義語が入力語となる場合は、その多義語に対応するすべての語義をしめすノードをネットワークに追加する。また、多義語同士間などいくつかの抑制リンクも追加する。文脈により関連のある語のネットワークが強化されていると考えられるため、3.2節で説明するノードの活性値の算出により、適切な語義ほど活性値が高くなる。

3.2 活性拡散を利用した語の多義性解消

3.1節で作成した Contextual Semantic Network のノードの活性値を、Interactive Activation Model (McClelland & Rumelhart, 1981) を参考にした式により算出する。ノードの活性値の最大値を1.0、最小値を0とする算出式を作成する。まずノードの初

期値 $a_i(0)$ を(3)式で表す。 S_{ki} は、入力文 k にノード N_i が出現した回数とする。また、時間 $t+1$ のときのノード N_i の値 ($a_i(t+1)$) は、(4)式とする。

$$a_i(0) = S_{ki} / 2, \tag{3}$$

$$a_i(t+1) = a_i(t) - \theta \cdot a_i(t) + \epsilon_i(t), \tag{4}$$

(4)式において、減衰係数 θ を (McClelland & Rumelhart, 1981) にならい0.1とし、 $\epsilon_i(t)$ は、時間 t において、隣接するノードから受ける影響を表す値を示す。ノード N_i に隣接するノードから受け取る活性値の合計 $n_i(t)$ は以下の(5)式のように定義する。

$$n_i(t) = \sum a_j(t) / \alpha D_{ij}, \tag{5}$$

ここで、 $a_j(t)$ はノード N_i に接続しているノード N_j の活性値を示し、 α は Contextual Semantic Network のリンク数の合計によって算出された補正值であり、 D_{ij} はノード N_i とノード N_j の間の距離としリンクの重みとして利用する。Contextual Semantic Network は、連想概念辞書の概念のネットワークで、

概念間の距離情報を持った連想関係のリンクを累積距離が 5.0 になるまで連想概念辞書の概念ネットワークを辿ることによって作成される。ネットワーク内で、抑制リンクとして働く場合はリンクの重み (D_{ij}) を負の値としている。

また、 $n_i(t) > 0$ で、入力全体が興奮状態になるとき、ノードへの影響 $\varepsilon_i(t)$ は、次の (6) 式で与えられる。

$$\varepsilon_i(t) = n_i(t) [M - a_i(t)], \quad (6)$$

ここで、 M は、ノードの最大活性化レベルであり、1.0 に設定している。

また、 $n_i(t) \leq 0$ で、入力全体が抑制状態にあるとき、ノードへの影響は次の (7) 式で与えられる。

$$\varepsilon_i(t) = n_i(t) [a_i(t) - m], \quad (7)$$

ここで、 m は、ノードの最小活性化レベルであり、0 に設定している。ノードの活性化値を計算し、多義語に相当するノード活性化値の大小関係で多義性の解消を行う。

3.3 連想概念辞書を利用したネットワークの構築

語の多義性解消を行うために、連想概念辞書の概念ネットワークを利用し 3.1 節の Contextual Semantic Network を以下のステップで作成する。

1. 多義語を含む入力文に対して、日本語係り受け解析器 Cabocha (Kudo & Matsumoto, 2002) を用いて、語の品詞情報と係り受け情報を取得し、文中の自立語を入力語 w_i とする。
2. 入力語 w_i が連想概念辞書の刺激語であれば、入力語 w_i の上位概念、部分材料、類義概念、環境概念としてリンクしている各概念を順に辿ることで語彙拡張を行ない、ネットワークに追加する。ただし、概念間の距離の累積距離が 5.0 になるまでとする。(入力語 w_i が多義語 A_i 、 B_i だとすると、 A_i と B_i を共に語彙拡張してネットワークに追加する。)
3. 入力語 w_i が連想概念辞書中の連想語であれば、

入力語 w_i とその刺激語 x_i までの距離を保存する。ステップ 2 と同様に刺激語 x_i から、リンクしている各概念を順に辿ることで語彙拡張を行ない、ネットワークに追加する。ただし、概念間の距離の累積距離が 5.0 になるまでとする。

4. 入力語 w_i の係り受け情報から、入力語 w_i の係り元の語とリンクをネットワークに追加し、リンクの重みを 1.0 とする。
(たとえば入力文が「絵を飾る」の時、「飾る」が入力語の場合は「飾る」の係り元である「絵」と「飾る」のリンクをネットワークに追加し、その重みを 1.0 とする。)
5. 多義語 (たとえば A_i と B_i) 同士のノードは抑制リンクをネットワークに追加し、リンクの重みを -1.0 とする。
6. 入力語 w_i が多義語 A_i 、 B_i だとすると A_i から連想される各々の連想語 a_j と B_i の間は抑制リンクをネットワークに追加し、リンクの重みを $-D_{ij}$ とする。

また同様に B_i から連想される各々の連想語 b_j と A_i の間は抑制リンクをネットワークに追加する。

図 3 の 2 つの長方形のノードは、多義語の「額」に値する 2 つの読み「がく」と「ひたい」を示す。多義語の場合は、複数の語義を示す刺激語が同時にネットワークに追加される。読みが「がく」の場合は「額縁」のみならず「金額」の意味の語もネットワークに追加される。これは、連想実験において「額(がく)」のように漢字とその読みを刺激語として呈示し、連想語に大きく影響を与えるようなインストラクションを与えていないためである。そのため連想概念辞書には「額縁」と「金額」に関する連想語が混在している。しかし本研究では同表記異音異義語でどちらの読みが適切であるかを判断するので、連想概念辞書はこの構造のまま利用する。

楕円内の語は「額(がく)」「額(ひたい)」を中心として、連想語(その連想語が刺激語として連想実験がなされている場合もある)の連想関係をたどりネットワークを作成した例である。実線は興奮リンクを示し、リンクの重みには距離情報 D_{ij} を用い

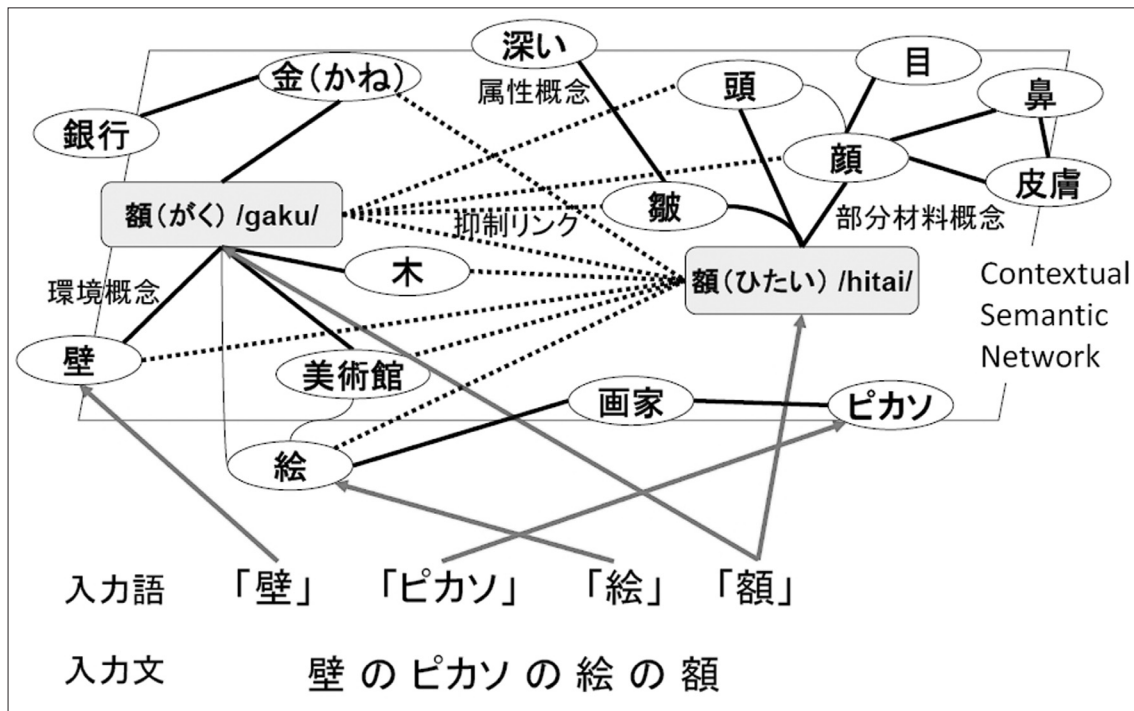


図3 連想概念辞書を用いた Contextual Semantic Network の例

る。点線は抑制リンクを示し、たとえば「額(がく)」と「額(ひたい)」の間のリンクの重みは-1.0としている。また、多義語 A_i から連想される a_j は多義語 B_i との抑制リンクを追加し、多義語 A_i から連想語 a_j までの距離情報に-1 を乗じた値をリンクの重み ($-D_{ij}$) としている。たとえば、図3の「木」と「額(ひたい) /hitai/」の点線は抑制リンクとなり、リンクの重みは「-1.46」となる。この場合のリンクの重みは、「木」と「額(がく) /gaku/」の距離に-1 を乗じた値である。

3.4 WordNet を利用したネットワークの構築

語の多義性解消を行うために、WordNet の概念構造を利用し3.1 節の Contextual Semantic Network を作成する。WordNet では synset という類語関係のセットで概念を定義している。たとえば、「額」を『ひたい』と読む場合は「frontal_bone、os_frontale、forehead」や「brow、forehead」の synset があり、前者は「the large cranial bone forming the front part of the cranium: includes the

upper part of the orbits」、後者は「the part of the face above the eyes」を指している。また、「額」を『がく』と読む場合は『金額』の意味で、「sum、sum_of_money、amount、amount_of_money」や「quantity」など複数の概念があり、前者は「a quantity of money」、後者は「an adequate or large amount」を指している。さらに『がく』と読み「額縁」の意味に対応する synset が WordNet にはないため、日本語の「額縁」に対応する WordNet の synset を利用する。

このように WordNet は日本語の表記に対応する synset が複数存在する。そこで WordNet を利用した Contextual Semantic Network で、多義語「額縁(がく)」、「金額(がく)」、「額(ひたい)」に対応する synset において、3.2 節の式を用いてノードの活性化値を計算し、多義語ごとの synset (ノード) の活性化値を平均し比較することで適切な語義を得るようなモデルを考案した。

表3は、WordNet の名詞において synset 間の関係性を示す記号とその意味を示す。「@ と ~」、「@i

表3 WordNet の名詞の関係子とその説明

関係子	説明	関係子	説明
@	Hypernym	+	Derivationally related form
@i	Instance Hypernym	;c	Domain of synset – TOPIC
~	Hyponym	-c	Member of this domain – TOPIC
~i	Instance Hyponym	;r	Domain of synset – REGION
#m	Member holonym	-r	Member of this domain – REGION
#s	Substance holonym	;u	Domain of synset – USAGE
#p	Part holonym	-u	Member of this domain – USAGE
%m	Member meronym	!	Antonym
%s	Substance meronym	=	Attribute
%p	Part meronym		

と~i]、[#mと%m]、[#sと%s]、[#pと%p]、[;cと-c]、[;rと-r]、[;uと-u]は双方向のリンクの関係にある。ここでは、語の多義性の解消を行うため、WordNetをもとにした Contextual Semantic Network を以下のステップで作成する。

1. 多義語を含む入力文に対して、日本語係り受け解析器 Cabocha を用いて、語の品詞情報と係り受け情報を取得し、文中の自立語を入力語 w_i とする。
2. 日本語 WordNet を利用して、入力語 w_i の日本語の表現に対応する WordNet の synset をすべて取得する。たとえば、「額」なら「額(ひたい)」に対応する synset として「frontal_bone、os_frontale、forehead」や「brow、forehead」など複数の概念がある。
3. 入力語 w_i に対応する synset の「@、@i、~、~i、#m、#s、#p、%m、%s、%p、+、;c、-c」の関係子を辿り関係子先の synset をノードとして得る。さらに累積経路が5になるまで順に「@」などの関係子先の synset をノードとして辿り語彙拡張を行ないネットワークに追加する。(入力語 w_i の synset が複数ある場合は、すべて利用して語彙拡張してネットワークに追加する。)
4. 入力語 w_i に関連性の高い synset を得るために、双方向のリンクを張るように「@と~」、「@iと~i」、「#mと%m」、「#sと%s」、「#pと%p」、「;cと-c」の組み合わせが経

- 路内にあるときには、同じ ID に戻ってくる経路をネットワークに追加する。たとえば、入力語が「鼻」の場合、synset である <nose、olfactory_organ> の part holonym <#p> に <face、human_face>、その hypernym <@> に <external_body_part>、その hypernym<@> に <body_part>、その part holonym <#p> に <organism、being>、その part meronym <%p> に <body_part> という経路が得られる。<body_part> がなんらかの synset (上記の場合は <organism、being>) を経由して再び <body_part> に戻るような経路である。
5. 入力語 w_i の係り受け情報から、入力語 w_i の係り元の語に対応する synset とリンクをネットワークに追加し、リンクの重みを 1.0 とする。
 6. また、多義語 A_i 、 B_i 、 C_i 同士の synset は抑制リンクをネットワークに追加し、リンクの重みを -1.0 とする。
たとえば、多義語「額(ひたい)」に対応する2つの synset のうちの 하나가「brow、forehead」の場合、他の語義を示す「額縁(がく)」、「金額(がく)」に対応する各 synset に抑制リンクで結ばれる。
 7. 入力語 w_i が多義語 A_i 、 B_i 、 C_i だとすると A_i のある一つの synset に隣接する synset a_j と多義語 B_i に対応する synset や多義語 C_i に対応する synset の間は抑制リンクをネットワークに追加し、リンクの重みを -1.0 とする。また同様に多義語 B_i に隣接する synset b_j と多義語 A_i に対応する synset や多

義語 C_i に対応する synset の間は抑制リンクをネットワークに追加する。

たとえば、多義語「額 (ひたい)」に対応する2つの synset のうちの 하나가「brow、forehead」の場合、上位概念である synset 「feature、lineament」は、他の語義である「額縁 (がく)」、「金額 (がく)」に対応する各 synset に抑制リンクで結ばれる。

図4は、WordNet を用いた Contextual Semantic Network の例である。実線は興奮リンクを示し、点線は抑制リンクを示す。また、丸は synset を示すノードである。長方形は、多義語 (「額縁 (がく)」、「金額 (がく)」、「額 (ひたい)」) に対応する複数の synset である。このように WordNet では日本語の表記に対して複数の synset の候補がある。たとえば、「額」なら『ひたい』の意味で「frontal_bone、os_frontale、forehead」や「brow、forehead」などの synset がある。連想概念辞書で「額 (ひたい)」に値するこの2つ synset は、図4の「額 (ひたい)

/hitai/」の長方形内のノードのように、synset 同士を興奮リンクで結ぶ。また、「額縁 (がく)」、「金額 (がく)」、「額 (ひたい)」内の各々の synset は、他の多義語同士と抑制リンクで結ばれる。たとえば、「brow、forehead」は「quantity」と抑制リンクで結ばれる。また、多義語に値する synset (「額縁 (がく)」、「金額 (がく)」、「額 (ひたい)」内のノード) から伸びた1段階の synset (たとえば、「brow、forehead」の場合「feature、lineament」を示すノード) は、他の語義を示す「額縁 (がく)」、「金額 (がく)」内の synset 一つ一つに抑制リンクで結ばれる。図4では、「feature、lineament」と「quantity」間の抑制リンクのことである。

3.5 連想概念辞書と WordNet を利用した多義性の解消システムの比較

3.3 節で、連想概念辞書の利用したネットワーク構造の作成について述べ、3.4 節で WordNet を利用したネットワーク構造の作成について述べた。多義語を含む分を入力文として、連想概念辞書を利

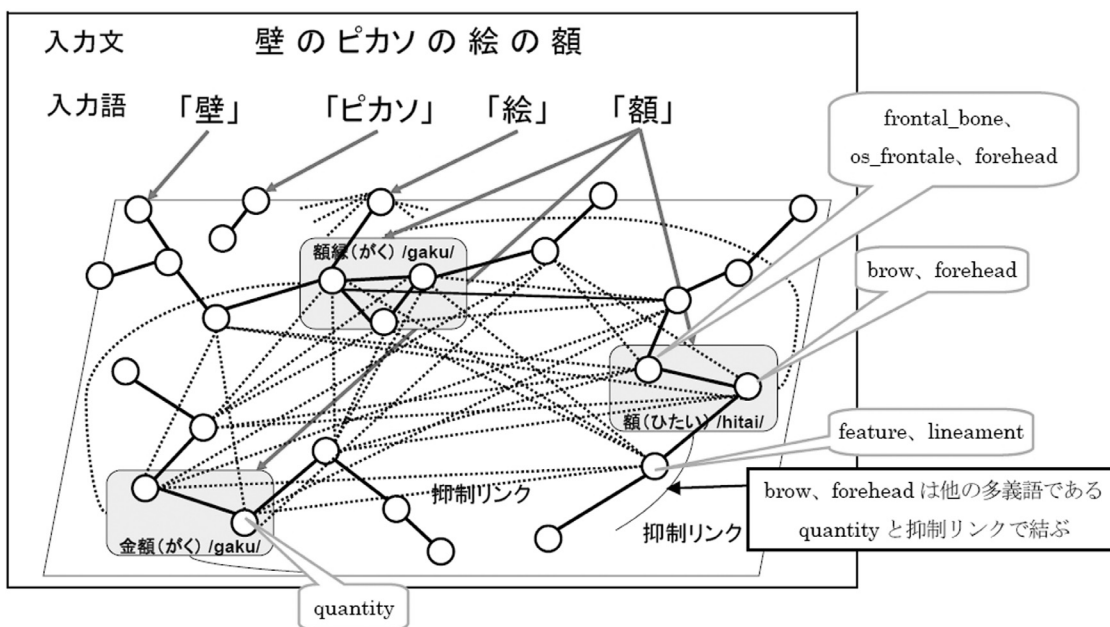


図4 WordNet を用いた Contextual Semantic Network の例 (抑制リンクは一部を省略)

用した時と WordNet を利用した時で Contextual Semantic Network を用いた多義性解消の正解率を比較する。ネットワーク中のノードの値は 3.2 節の式を利用して計算する。

3.5.1 多義性解消手法の比較の手順

Contextual Semantic Network を利用した多義性の解消の評価を行うために、多義語「額」が文中に 1 語存在する Web 上の文約 2 万 8 千を取得し、そのうち連想概念辞書の刺激語にある概念で文の名詞が構成されている 26 文を選定した。額（がく）と読む文は 16 文で、額（ひたい）と読む文は 10 文である。額（がく）と読む文については 16 文のうち、「額縁」の意味の文は 12 文で、残りの 4 文は、「金額」の意味の文である。選定した 26 文を入力文として、文ごとに Contextual Semantic Network を作成したうえで多義語のノードの活性値を算出し、ノードの活性値が高い方をシステムが選択した語義とする。Contextual Semantic Network は連想概念辞書を利用した場合と WordNet を利用した場合で作成し、システムが選択した語義が本来の適切な語義であるか正解率を比較する。また、入力文ごとに連想概念辞書と WordNet でネットワーク密度を計算した。

3.5.2 結果と考察

表 4 は、3.5.1 項で選定した 26 文を利用して、連想概念辞書と WordNet を利用して作成した Contextual Semantic Network をもとにした多義性解消の正解率を示す。連想概念辞書は「額（がく）」の場合、被験者に提示した刺激語が「額（がく）」となるため、「金額」に関する語も「額縁」に関する語も連想されており 2 つは明確に分類されていないが、WordNet は synset として概念が明確に分類

されている。本論文では、同表記異音異義語を対象にした語の多義性解消を行っており、正しい読みが判断できるかを目的としているため、連想概念辞書のデータをそのまま利用している。

ネットワーク密度を比較すると、連想概念辞書は WordNet よりもネットワーク密度が 10 倍ほど高い。2.3.2 項の表 2 における刺激語のみを対象とした連想概念辞書のネットワーク密度よりも低くなっている。これは 3.3 節のステップで作成する Contextual Semantic Network が連想語も含めたネットワークとなっているためである。しかし刺激語のみを対象とした連想概念辞書のネットワーク密度よりも低いとはいえ、連想概念辞書の方が WordNet よりも多くの概念が密につながっていることが分かる。

正解率は、連想概念辞書の場合が 92.3% で WordNet の場合が 58.3% となり、連想概念辞書を用いた多義性の解消の正解率の方が WordNet を用いた多義性の解消の正解率より高い。連想概念辞書には、「額（がく）」の環境概念は「壁」など、「額（がく）」がどのような状況・環境にあるかなどの文脈情報なども記載されている。WordNet にも「;c Domain of synset – TOPIC」のようにどのような Domain で用いられる概念であるかなどの情報もあるが、網羅数が少ない。そこで、ネットワーク密度を上げるには、多量のコーパスなどを利用して、「壁にかかる絵」「壁の絵」など、係り受け情報や表現のパターンなどを活用して語の共起情報から、文脈情報を抽出し、WordNet を補完する方法が考えられる。

4 動的文脈ネットワークモデル (Dynamic Contextual Network Model) と多義性の解消

2 章では、連想概念辞書の構築について述べ、さ

表 4 連想概念辞書と WordNet を利用した多義性の解消精度の比較

漢字的読み	連想概念辞書		WordNet	
	額（がく）	額（ひたい）	額（がく・額縁） 額（がく・金額）	額（ひたい）
平均ネットワーク密度	0.0178		0.0018	
正解数	15	9	8	6
正解率	92.3%		58.3%	

らに連想概念辞書と WordNet の規模とネットワーク密度の比較を行った。3章では、文全体で文脈ネットワークを作成して連想概念辞書と WordNet で多義性解消の正解率を比較した。次に本章では、動的文脈ネットワークモデル (Okamoto, et. al., 2010) (Okamoto, et. al., 2011) を提案する。このモデルは、人間が文を読み進めながら内容を理解し多義性を解消するように、文中の単語を入力順に用いて、それぞれの時点での文脈に応じたネットワーク (Contextual Semantic Network) を動的に構築し、曖昧性を随時解消するモデルである。動的文脈ネットワークモデルは、ネットワークを作ったある時点で語義の同定が十分でない場合、さらに次の入力語を加えネットワークを再構築しながら、活性値を計算することで語の多義性を解消する。文中の単語を順番に利用してネットワークを動的に作成するためネットワーク密度が高いことが求められる。そのため連想概念辞書のみを利用し語の多義性解消を行う。

4.1 ネットワークの動的変化

図5は、動的文脈ネットワークモデルの特徴を示す。入力単語列を確定した時点での Contextual Semantic Network を作成し多義性の解消を試みた結果、どの語義が適切か判断できなかった場合、次の入力語を用いて語彙拡張を行う。こうして構築したネットワークで活性値を再計算することにより適切な語義を選択する。

図5の左側部分は、「額に入れて絵を飾る」の

ように、文の先頭部分が多義語である場合の Contextual Semantic Network をどのように拡張するか示している。最初の「額」だけではどちらの読みが適切であるか判断できないため、次の「入れる」や「絵」といった入力語を用いて拡張したネットワークで随時ノードの活性値を計算することで適切な語義を判断する。具体的には、語が入力されるたびに3.2節の式を利用して、時間 t については語が入力されるたびに20サイクル計算する。ネットワーク内の最初のノードの初期値は3.2節の(3)式で算出するが、次の語が入力された場合のノードの初期値は、その時点の活性値の値を継承して(8)式で計算する

$$a_i(0) = S_{ki} / 2, \tag{3}$$

$$a_i(t) = (S_{ki} + a_i(t-1)) / 2, \tag{8}$$

S_{ki} は、文 k にノード N_i が出現した回数である。 $a_i(t-1)$ は、時間 $t-1$ におけるノード N_i の活性値である。また、時間 $t+1$ におけるノード N_i の値 $a_i(t+1)$ は、3.2節の(4)から(7)式を用いて計算する。

4.2 動的文脈ネットワークモデルを利用した多義性解消のシミュレーション

4.2.1 多義性解消のシミュレーションと評価方法

動的文脈ネットワークモデルを利用すると、「壁にかかったピカソの絵の額が、頭に当たって額から血が出た。」のように、一つの文の中に多義語が複数ある場合や、「額ににじむ汗を手で拭う。写真は

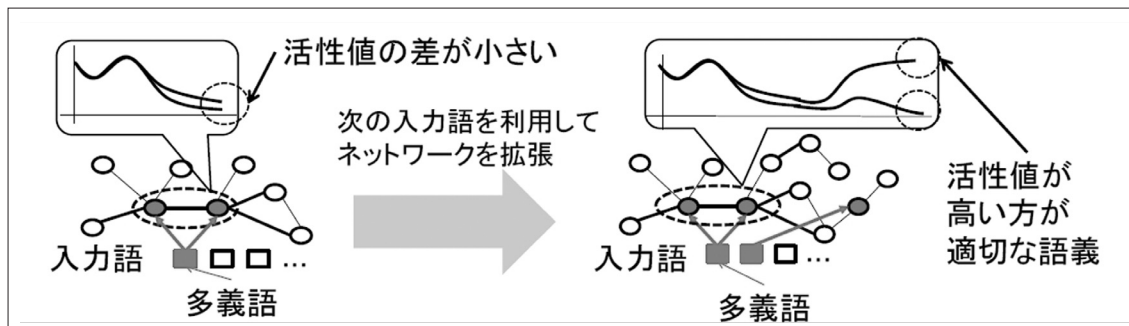


図5 Contextual Semantic Network の拡張

額に入れられ壁に飾られた。」のように2つの文を組み合わせた場合や、冒頭が多義語「額」で正しい読みを判断できない場合でも、後の文脈を取り入れることによって多義性の解消に対応できる。入力文が複数の文で構成される場合は、文の切れ目に関係なくノードの活性値は次の文のネットワークにも継承される。

図6は、「額ににじむ汗を手で拭う。写真は額に入れられ壁に飾られた。」を入力文として、動的文脈ネットワークモデルを利用した語の多義性解消のシミュレーション結果である。最初に「額（がく）」と「額（ひたい）」が同時に入力され「汗」が入力されるまで「額（がく）」と「額（ひたい）」の活性値の差がない。その後は、丸い点線の箇所のように「額（ひたい）」が「額（がく）」より高く、この時点では「額（ひたい）」と読むのが適切であると読み取れる。丸い実線の箇所も「額（がく）」、「額（ひたい）」が同時に入力された部分である。「額（がく）」が「額（ひたい）」より高く、この時点では「額（がく）」と読むのが適切であると読み取れる。動的文脈ネットワークモデルの評価として3章で用いた26文（「がく」と読む16文と「ひたい」と読む10文）の「額」を用い、「がく」と読む文と「ひたい」と読む文の2つを組み合わせた図6の「額に

にじむ汗を手で拭う。写真は額に入れられ壁に飾られた。」のような入力文を320個用意した。また「額」と同様、多義語「金」が文中に1語存在するWeb上で取得された約6万8千文中、連想概念辞書の刺激語にある名詞で構成されている文を22文（「きん」と読む12文と「かね」と読む10文）用い、2つの文を組み合わせた240個の入力文を用意した。前の文と後の文に含まれる多義語について、動的文脈ネットワークモデルを用いて適切な語義が得られるかを調べた。

4.2.2 結果と考察

表5は、多義語を含む2つの文をつなげて、動的文脈ネットワークモデルで解析した結果である。入力文中の2つの文の2つの多義語を両方とも正しく解析できた場合を「正解数2」とし、2つの文のうち片方の多義語しか正しく解析できなかった場合を「正解数1」とし、2文とも多義語を正しく解析できなかった場合を「正解数0」としている。また、入力文に含まれる2か所の多義語の語義（「額」の場合は320個×2=640か所、「金」の場合は240個×2=480か所）の総数に対して、語義を正しく判断できた割合を正解率とした。

「額」と「金」において多義を含む文が「正解数2」

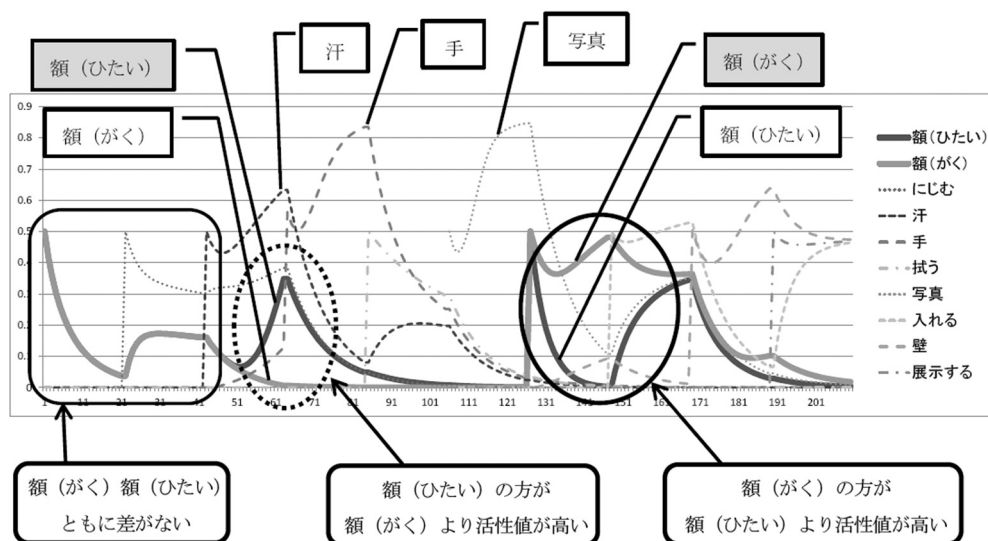


図6 「額ににじむ汗を手で拭う。写真は額に入れられ壁に飾られた。」を入力文としたときの語の活性値の変化

表5 多義語を含む2つの文を組み合わせた場合の提案モデルにおける正解数

多義語	入力文	正解数 2	正解数 1	正解数 0	正解率
額 (がく) / 額 (ひたい)	320 個	222 個 (69.6%)	92 個 (28.6%)	6 個 (1.8%)	83.8%
金 (きん) / 金 (かね)	240 個	141 個 (58.8%)	99 個 (41.2%)	0 個 (0%)	79.4%

だったのは約6割から7割であった。また、入力文の前半の文もしくは後半の文の「正解数1」だったのは約3割から4割であった。

第3章で述べたように、1文中の語すべてで Contextual Semantic Network を作成して多義性の解消をした場合は正解率が9割と高かった。語義が異なる多義語を含む2つの文を組み合わせて作成した入力文に対して、動的に Contextual Semantic Network を作り変えながらノードの活性値を計算した場合に正解率は8割前後と劣るものの、複数の多義語が存在する場合であっても前後の文脈を得ながら多義性の解消ができることが分かる。「正解数1」で誤りが顕著だったのは、後半の文の文頭に多義語が来る場合である。前の文のノードの活性値を継承して2文目を読み込むため文頭に多義語があった場合に誤った判断をすること多かった。今後は、多義語の位置より後ろの文脈によっては、一度決定した判断を取り消して適切な多義語を選択する仕組みを取り入れる必要があると考える。また、一語ずつ入力語として利用すると入力数が少なく文脈情報を十分にネットワークに反映することができないので、Cabochaの係り受け情報を利用し句や節単位でネットワークを再構築することも進めていく必要がある。

5 まとめと今後の展開

本論文において、2章では連想概念辞書の特徴と構築方法について述べ、規模やネットワーク密度を WordNet と比較した。3章では多義語を含む文全体を対象にした文脈ネットワークモデルによる多義性解消法を提案した。それは、入力文から Contextual Semantic Network を作成し、活性拡散によって語の多義性を解消する方法である。Contextual Semantic Network は連想概念辞書を利用した場合と WordNet を利用した場合とで多義性解消の正解

率を比較し、連想概念辞書を利用した方が高い正解率を得た。4章では、人間が入力文の文頭から順番に語を見ながら意味を理解していくように、入力される語を順番に文脈ネットワークに組み入れて再構築する動的な文脈ネットワークモデルを提案した。また曖昧性のある語の入力時点で、多義語の意味が判断できない場合は、さらに次の語もネットワークに取り込んでから語義を理解するように設計した。同表記異音異義語を含む2文を組み合わせて入力文とし、同モデルを利用した多義性の解消において約8割の正解率を得た。

連想概念辞書は、概念間の距離情報を持つと同時に、既存概念辞書にはない環境概念(状況)や属性概念、動作概念などの連想関係にある語も体系化した辞書である。今後は、WordNetなどの他の概念辞書と連想概念辞書とを組み合わせた辞書を用いて文脈ネットワークを作成し、連想概念辞書を追加することの有効性について調査を行っていきたい。

参考文献

- 1 A. M. Collins and M. R. Quillian, "Retrieval Time from Semantic Memory," *Journal of Verbal Learning and Verbal Behavior*, vol. 8, pp.240-247, 1969.
- 2 A. M. Collins and E. F. Loftus, "A Spreading-Activation Theory of Semantic Processing," *Psychological Review*, vol.82-6, pp.407-428, 1975.
- 3 EDR, 『電子化辞書使用説明書』, 1990年.
- 4 H. Isahara, F. Bond, K. Uchimoto, M. Utiyama and K. Kanzaki, "Development of Japanese WordNet," *Proc of LREC 2008*, 2008.
- 5 甲斐 睦朗・松川 利広, 『語彙指導の方法』, 光村図書, 1995年.
- 6 T. Kudo and Y. Matsumoto, "Japanese Dependency Analysis using Cascaded Chunking," *Proc. of CONLL 2002*, pp.63-69, 2002.
- 7 J. L. McClelland and D. E. Rumelhart, "An Interactive Activation Model of Context Effects in Letter Perception: Part 1. An Account of Basic Findings," *Psychological Rev.*, Vol.88, No.5, pp.375-407, 1981.
- 8 G. A. Miller, R. Beckwin, C. Fellbaum and D. Gross, K. Miller and R. Teng, "Five Papers on WordNet," *CSL Report 43*, Cognitive Science Laboratory Princeton

- University, 1993.
- 9 村田 真樹・内山 将夫・内元 清貴・馬青・井佐原 均, 「SENSEVAL2」辞書タスクのCRLの取り組みー日本語単語の多義性解消における種々の機械学習手法と素性の比較ー, 『自然言語処理』, Vol.10, No.3, pp.115-133, 2003年.
 - 10 岡本 潤・石崎 俊, 「概念間距離の定式化と電子化辞書との比較」, 『自然言語処理』, Vol.8, No.4, pp.37-54, 2001年.
 - 11 岡本 潤・石崎 俊, 「連想概念辞書の距離情報を用いた重要文の抽出」, 『自然言語処理』, Vol.10, No.5, pp.139-151, 2003年.
 - 12 J. Okamoto and S. Ishizaki, "Homographic Ideogram Understanding Using Contextual Dynamic Network," *Proc. of LREC 2010*, 2010.
 - 13 J. Okamoto and S. Ishizaki, "An Associative Concept Dictionary for Natural Language Processing: Text Summarization and Word Sense Disambiguation," *Journal of Cognitive Science*, Vol. 12, pp. 259-276, 2011.
 - 14 E. Rosch and C. B. Mervis, "Family Resemblances: Studies in the Internal Structure of Categories," *Cognitive Psychology*, Vol.7, pp.573-605, 1975.
 - 15 坂口 琢哉, 「連想概念辞書のニューラルネットワークへの符号化と比喩理解システムの応用」, 『安田女子大学紀要』, No.38, pp.169-179, 2010年.
 - 16 J. Veronis and N. M. Ide, "Word Sense Disambiguation with Very Large Neural Networks Extracted from Machine Readable Dictionaries," *Proc. of Coling 1990*, pp. 389-394, 1990.
 - 17 D. L. Waltz and J. B. Pollack, "Massively parallel parsing: A strongly interactive model of natural language interpretation," *Cognitive Science*, Vol. 9, pp. 51-74, 1985.

[2012. 3. 30 受理]

[2012. 7. 6 採録]